

## Class Notes Econ 7800 Fall Semester 2003

Hans G. Ehrbar

ECONOMICS DEPARTMENT, UNIVERSITY OF UTAH, 1645 CAMPUS CENTER  
DRIVE, SALT LAKE CITY UT 84112-9300, U.S.A.

*URL:* [www.econ.utah.edu/ehrbar/ecmet.pdf](http://www.econ.utah.edu/ehrbar/ecmet.pdf)

*E-mail address:* [ehrbar@econ.utah.edu](mailto:ehrbar@econ.utah.edu)

ABSTRACT. This is an attempt to make a carefully argued set of class notes freely available. The source code for these notes can be downloaded from [www.econ.utah.edu/ehrbar/ecmet-sources.zip](http://www.econ.utah.edu/ehrbar/ecmet-sources.zip) Copyright Hans G. Ehrbar under the GNU Public License

The present version has those chapters relevant for Econ 7800.

## Contents

Chapter 1. Syllabus Econ 7800 Fall 2003	ix	Chapter 5. Chebyshev Inequality, Weak Law of Large Numbers, and Central Limit Theorem	
Chapter 2. Probability Fields	1	5.1. Chebyshev Inequality	
2.1. The Concept of Probability	1	5.2. The Probability Limit and the Law of Large Numbers	
2.2. Events as Sets	6	5.3. Central Limit Theorem	
2.3. The Axioms of Probability	10	Chapter 6. Vector Random Variables	
2.4. Objective and Subjective Interpretation of Probability	13	6.1. Expected Value, Variances, Covariances	
2.5. Counting Rules	14	6.2. Marginal Probability Laws	
2.6. Relationships Involving Binomial Coefficients	16	6.3. Conditional Probability Distribution and Conditional Mean	
2.7. Conditional Probability	17	6.4. The Multinomial Distribution	
2.8. Ratio of Probabilities as Strength of Evidence	23	6.5. Independent Random Vectors	
2.9. Bayes Theorem	24	6.6. Conditional Expectation and Variance	
2.10. Independence of Events	25	6.7. Expected Values as Predictors	
2.11. How to Plot Frequency Vectors and Probability Vectors	28	6.8. Transformation of Vector Random Variables	
Chapter 3. Random Variables	33	Chapter 7. The Multivariate Normal Probability Distribution	
3.1. Notation	33	7.1. More About the Univariate Case	
3.2. Digression about Infinitesimals	33	7.2. Definition of Multivariate Normal	
3.3. Definition of a Random Variable	35	7.3. Special Case: Bivariate Normal	
3.4. Characterization of Random Variables	36	7.4. Multivariate Standard Normal in Higher Dimensions	
3.5. Discrete and Absolutely Continuous Probability Measures	40	Chapter 8. The Regression Fallacy	
3.6. Transformation of a Scalar Density Function	41	Chapter 9. A Simple Example of Estimation	
3.7. Example: Binomial Variable	43	9.1. Sample Mean as Estimator of the Location Parameter	
3.8. Pitfalls of Data Reduction: The Ecological Fallacy	44	9.2. Intuition of the Maximum Likelihood Estimator	
3.9. Independence of Random Variables	45	9.3. Variance Estimation and Degrees of Freedom	
3.10. Location Parameters and Dispersion Parameters of a Random Variable	46	Chapter 10. Estimation Principles and Classification of Estimators	
3.11. Entropy	52	10.1. Asymptotic or Large-Sample Properties of Estimators	
Chapter 4. Specific Random Variables	63	10.2. Small Sample Properties	
4.1. Binomial	63		
4.2. The Hypergeometric Probability Distribution	66		
		4.3. The Poisson Distribution	
		4.4. The Exponential Distribution	
		4.5. The Gamma Distribution	
		4.6. The Uniform Distribution	
		4.7. The Beta Distribution	
		4.8. The Normal Distribution	
		4.9. The Chi-Square Distribution	
		4.10. The Lognormal Distribution	
		4.11. The Cauchy Distribution	

10.3. Comparison Unbiasedness Consistency	154
10.4. The Cramer-Rao Lower Bound	158
10.5. Best Linear Unbiased Without Distribution Assumptions	166
10.6. Maximum Likelihood Estimation	168
10.7. Method of Moments Estimators	171
10.8. M-Estimators	171
10.9. Sufficient Statistics and Estimation	171
10.10. The Likelihood Principle	175
10.11. Bayesian Inference	176
Chapter 11. Interval Estimation	179
Chapter 12. Hypothesis Testing	187
12.1. Duality between Significance Tests and Confidence Regions	190
12.2. The Neyman Pearson Lemma and Likelihood Ratio Tests	191
12.3. The Wald, Likelihood Ratio, and Lagrange Multiplier Tests	194
Chapter 13. General Principles of Econometric Modelling	197
Chapter 14. Mean-Variance Analysis in the Linear Model	199
14.1. Three Versions of the Linear Model	199
14.2. Ordinary Least Squares	200
14.3. The Coefficient of Determination	208
14.4. The Adjusted R-Square	213
Chapter 15. Digression about Correlation Coefficients	215
15.1. A Unified Definition of Correlation Coefficients	215
Chapter 16. Specific Datasets	219
16.1. Cobb Douglas Aggregate Production Function	219
16.2. Houthakker's Data	229
16.3. Long Term Data about US Economy	234
16.4. Dougherty Data	235
16.5. Wage Data	236
Chapter 17. The Mean Squared Error as an Initial Criterion of Precision	253
17.1. Comparison of Two Vector Estimators	253
Chapter 18. Sampling Properties of the Least Squares Estimator	257
18.1. The Gauss Markov Theorem	258
18.2. Digression about Minimax Estimators	260
18.3. Miscellaneous Properties of the BLUE	261
18.4. Estimation of the Variance	271

18.5. Mallow's Cp-Statistic as Estimator of the Mean Squared Error	272
Chapter 19. Nonspherical Positive Definite Covariance Matrix	272
Chapter 20. Best Linear Prediction	272
20.1. Minimum Mean Squared Error, Unbiasedness Not Required	272
20.2. The Associated Least Squares Problem	272
20.3. Prediction of Future Observations in the Regression Model	272
Chapter 21. Updating of Estimates When More Observations become Available	272
Chapter 22. Constrained Least Squares	272
22.1. Building the Constraint into the Model	272
22.2. Conversion of an Arbitrary Constraint into a Zero Constraint	272
22.3. Lagrange Approach to Constrained Least Squares	272
22.4. Constrained Least Squares as the Nesting of Two Simpler Models	272
22.5. Solution by Quadratic Decomposition	272
22.6. Sampling Properties of Constrained Least Squares	272
22.7. Estimation of the Variance in Constrained OLS	272
22.8. Inequality Restrictions	272
22.9. Application: Biased Estimators and Pre-Test Estimators	272
Chapter 23. Additional Regressors	272
Chapter 24. Residuals: Standardized, Predictive, "Studentized"	272
24.1. Three Decisions about Plotting Residuals	272
24.2. Relationship between Ordinary and Predictive Residuals	272
24.3. Standardization	272
Chapter 25. Regression Diagnostics	272
25.1. Missing Observations	272
25.2. Grouped Data	272
25.3. Influential Observations and Outliers	272
25.4. Sensitivity of Estimates to Omission of One Observation	272
Chapter 26. Asymptotic Properties of the OLS Estimator	272
26.1. Consistency of the OLS estimator	272
26.2. Asymptotic Normality of the Least Squares Estimator	272
Chapter 27. Least Squares as the Normal Maximum Likelihood Estimate	272
Chapter 28. Random Regressors	272

28.1. Strongest Assumption: Error Term Well Behaved Conditionally on Explanatory Variables	355
28.2. Contemporaneously Uncorrelated Disturbances	357
28.3. Disturbances Correlated with Regressors in Same Observation	357
Chapter 29. The Mahalanobis Distance	359
29.1. Definition of the Mahalanobis Distance	359
Chapter 30. Interval Estimation	363
30.1. A Basic Construction Principle for Confidence Regions	363
30.2. Coverage Probability of the Confidence Regions	367
30.3. Conventional Formulas for the Test Statistics	368
30.4. Interpretation in terms of Studentized Mahalanobis Distance	368
Chapter 31. Three Principles for Testing a Linear Constraint	373
31.1. Mathematical Detail of the Three Approaches	374
31.2. Examples of Tests of Linear Hypotheses	377
31.3. The F-Test Statistic is a Function of the Likelihood Ratio	385
31.4. Tests of Nonlinear Hypotheses	386
31.5. Choosing Between Nonnested Models	386
Chapter 32. Instrumental Variables	387
Appendix A. Matrix Formulas	393
A.1. A Fundamental Matrix Decomposition	393
A.2. The Spectral Norm of a Matrix	394
A.3. Inverses and g-Inverses of Matrices	394
A.4. Deficiency Matrices	396
A.5. Nonnegative Definite Symmetric Matrices	399
A.6. Projection Matrices	404
A.7. Determinants	406
A.8. More About Inverses	407
A.9. Eigenvalues and Singular Value Decomposition	411
Appendix B. Arrays of Higher Rank	415
B.1. Informal Survey of the Notation	415
B.2. Axiomatic Development of Array Operations	418
B.3. An Additional Notational Detail	423
B.4. Equality of Arrays and Extended Substitution	423
B.5. Vectorization and Kronecker Product	424
Appendix C. Matrix Differentiation	435

C.1. First Derivatives

Appendix. Bibliography

## CHAPTER 1

## Syllabus Econ 7800 Fall 2003

The class meets Tuesdays and Thursdays 12:25 to 1:45pm in BUC 207. First class Thursday, August 21, 2003; last class Thursday, December 4.

Instructor: Assoc. Prof. Dr. Hans G. Ehrbar. Hans's office is at 319 BUO, Tel. 581 7797, email [ehrbar@econ.utah.edu](mailto:ehrbar@econ.utah.edu) *Office hours:* Monday 10–10:45 am, Thursday 5–5:45 pm or by appointment.

*Textbook:* There is no obligatory textbook in the Fall Quarter, but detailed class notes are available at [www.econ.utah.edu/ehrbar/ec7800.pdf](http://www.econ.utah.edu/ehrbar/ec7800.pdf), and you can purchase a hardcopy containing the assigned chapters only at the University Copy Center, 158 Union Bldg, tel. 581 8569 (ask for the class materials for Econ 7800).

Furthermore, the following optional texts will be available at the bookstore: Peter Kennedy, *A Guide to Econometrics* (fourth edition), MIT Press, 1998 ISBN 0-262-61140-6.

The bookstore also has available William H. Greene's *Econometric Analysis*, fifth edition, Prentice Hall 2003, ISBN 0-13-066189-9. This is the assigned text for Econ 7801 in the Spring semester 2004, and some of the introductory chapters are already useful for the Fall semester 2003.

The following chapters in the class notes are assigned: 2, 3 (but not section 3.2), 4, 5, 6, 7 (but only until section 7.3), 8, 9, 10, 11, 12, 14, only section 15.1 in chapter 15, in chapter 16, we will perhaps do section 16.1 or 16.4, then in chapter 17 we do section 17.1, then chapter 18 until and including 18.5, and in chapter 22 do sections 22.1, 22.3, 22.6, and 22.7. In chapter 29 only the first section 29.1, finally chapters 30, and section 31.2 in chapter 31.

*Summary of the Class:* This is the first semester in a two-semester Econometrics field, but it should also be useful for students taking the first semester only as part of their methodology requirement. The course description says: Probability, conditional probability, distributions, transformation of probability densities, sufficient statistics, limit theorems, estimation principles, maximum likelihood estimation, interval estimation and hypothesis testing, least squares estimation, linear constraints.

This class has two focal points: maximum likelihood estimation, and the fundamental concepts of the linear model (regression).

If advanced mathematical concepts are necessary in these theoretical explanations, they will usually be reviewed very briefly before we use them. The class is structured in such a way that, if you allocate enough time, it should be possible to refresh your math skills as you go along.

Here is an overview of the topics to be covered in the Fall Semester. They may not come exactly in the order in which they are listed here

**1. Probability fields:** Events as sets, set operations, probability axioms, subjective vs. frequentist interpretation, finite sample spaces and counting rules (combinatorics), conditional probability, Bayes theorem, independence, conditional independence.

**2. Random Variables:** Cumulative distribution function, density function, location parameters (expected value, median) and dispersion parameters (variance)

**3. Special Issues and Examples:** Discussion of the “ecological fallacy”; entropy; moment generating function; examples (Binomial, Poisson, Gamma, Normal, Chisquare); sufficient statistics.

**4. Limit Theorems:** Chebyshev inequality; law of large numbers; central limit theorems.

The first Midterm will already be on Thursday, September 18, 2003. It will be a closed book, but you are allowed to prepare one sheet with formulas etc. Most of the midterm questions will be similar or identical to the homework questions in the class notes assigned up to that time.

**5. Jointly Distributed Random Variables:** Joint, marginal, and conditional densities; conditional mean; transformations of random variables; covariance and correlation; sums and linear combinations of random variables; jointly normal variables.

**6. Estimation Basics:** Descriptive statistics; sample mean and variance; degrees of freedom; classification of estimators.

**7. Estimation Methods:** Method of moments estimators; least squares estimators. Bayesian inference. Maximum likelihood estimators; large sample properties of MLE; MLE and sufficient statistics; computational aspects of maximum likelihood estimation.

**8. Confidence Intervals and Hypothesis Testing:** Power functions; Neyman Pearson Lemma; likelihood ratio tests. As example of tests: the run test, goodness of fit test, contingency tables.

The second in-class Midterm will be on Thursday, October 16, 2003.

**9. Basics of the “Linear Model.”** We will discuss the case with nonrandom regressors and a spherical covariance matrix: OLS-BLUE duality, Maximum likelihood estimation, linear constraints, hypothesis testing, interval estimation (t-test, F-test, joint confidence intervals).

The third Midterm will be a takehome exam. You will receive the questions on Tuesday, November 25, 2003, and they are due back at the beginning of class

Tuesday, December 2nd, 12:25 pm. The questions will be similar to questions which you might have to answer in the Econometrics Field exam.

The Final Exam will be given according to the campus-wide examination schedule, which is Wednesday December 10, 10:30–12:30 in the usual classroom. Closed book, but again you are allowed to prepare one sheet of notes with the most important concepts and formulas. The exam will cover material after the second Midterm.

*Grading:* The three midterms and the final exams will be counted equally. Every week certain homework questions from among the questions in the class notes will be assigned. It is recommended that you work through these homework questions conscientiously. The answers provided in the class notes should help you if you get stuck. If you have problems with these homeworks despite the answers in the class notes, please write your answer down as far as you get and submit your answer to me; I will look at them and help you out. A majority of the questions in the two in-class midterms and the final exam will be identical to these assigned homework questions, but some questions will be different.

*Special circumstances:* If there are special circumstances requiring an individualized course of study in your case, please see me about it in the first week of classes.

Hans G. Ehrbar

## CHAPTER 2

## Probability Fields

## 2.1. The Concept of Probability

Probability theory and statistics are useful in dealing with the following types of situations:

- Games of chance: throwing dice, shuffling cards, drawing balls out of urns.
- Quality control in production: you take a sample from a shipment, count how many defectives.
- Actuarial Problems: the length of life anticipated for a person who has just applied for life insurance.
- Scientific Experiments: you count the number of mice which contract cancer when a group of mice is exposed to cigarette smoke.
- Markets: the total personal income in New York State in a given month.
- Meteorology: the rainfall in a given month.
- Uncertainty: the exact date of Noah's birth.
- Indeterminacy: The closing of the Dow Jones industrial average or the temperature in New York City at 4 pm. on February 28, 2014.
- Chaotic determinacy: the relative frequency of the digit 3 in the decimal representation of  $\pi$ .
- Quantum mechanics: the proportion of photons absorbed by a polarization filter
- Statistical mechanics: the velocity distribution of molecules in a gas at a given pressure and temperature.

In the probability theoretical literature the situations in which probability theory applies are called “experiments,” see for instance [Rén70, p. 1]. We will not use this terminology here, since probabilistic reasoning applies to several different types of situations, and not all these can be considered “experiments.”

PROBLEM 1. (*This question will not be asked on any exams*) Rényi says: “Observing how long one has to wait for the departure of an airplane is an experiment.”  
Comment.

ANSWER. Rényi commits the epistemic fallacy in order to justify his use of the word “experiment.” Not the observation of the departure but the departure itself is the event which can be theorized probabilistically, and the word “experiment” is not appropriate here.

What does the fact that probability theory is appropriate in the above situations tell us about the world? Let us go through our list one by one:

- Games of chance: Games of chance are based on the sensitivity on initial conditions: you tell someone to roll a pair of dice or shuffle a deck of cards and despite the fact that this person is doing exactly what he or she is asked to do and produces an outcome which lies within a well-defined universe known beforehand (a number between 1 and 6, or a permutation of the deck of cards), the question *which* number or *which* permutation is beyond their control. The precise location and speed of the die or the precise order of the cards varies, and these small variations in initial conditions give rise by the “butterfly effect” of chaos theory, to unpredictable final outcomes.

A critical realist recognizes here the openness and stratification of the world: If many different influences come together, each of which is governed by laws, then their sum total is not determinate, as a naive hyper-determinist would think, but indeterminate. This is not only a condition for the possibility of science (in a hyper-deterministic world, one could not know anything before one knew everything, and science would also not be necessary because one could not do anything), but also for practical human activity: the macro outcomes of human practice are largely independent of micro detail (the postcard arrives whether the address is written in cursive or in printed letters, etc.). Games of chance are situations which deliberately project this micro indeterminacy into the macro world: the micro influences cancel each other out without one enduring influence taking over (as would be the case if the die were not perfectly symmetric and balanced) or deliberate human corrective activity stepping into the void (as a card-trickster might do if the cards being shuffled somehow were distinguishable from the backside).

The experiment in which one draws balls from urns shows clearly another aspect of this paradigm: the set of different possible outcomes is fixed beforehand, and the probability enters in the choice of one of the predetermined outcomes. This is not the only way probability can arise: it is an extensionalist example, in which the connection between success and failure is external. The world is not a collection of externally related outcomes collected in an urn. Success and failure are not determined by choice between different spatially separated and individually inert balls (like playing cards or faces on a die), but it is the outcome of development and struggle that is internal to the individual unit.

- Quality control in production: you take a sample from a shipment, count how many defectives. Why is statistics and probability useful in production? Because production is work, it is not spontaneous. Nature does not voluntarily give us things in the form in which we need them. Production is similar to a scientific experiment because it is the attempt to create local closure. Such closure can never be complete, there are always leaks in it, through which irregularity enters.
- Actuarial Problems: the length of life anticipated for a person who has just applied for life insurance. Not only production, but also life itself is a struggle with physical nature, it is emergence. And sometimes it fails: sometimes the living organism is overwhelmed by the forces which it tries to keep at bay and to subject to its own purposes.
- Scientific Experiments: you count the number of mice which contract cancer when a group of mice is exposed to cigarette smoke: There is local closure regarding the conditions under which the mice live, but even if this closure were complete, individual mice would still react differently, because of genetic differences. No two mice are exactly the same, and despite these differences they are still mice. This is again the stratification of reality. Two mice are two different individuals but they are both mice. Their reaction to the smoke is not identical, since they are different individuals, but it is not completely capricious either, since both are mice. It can be predicted probabilistically. Those mechanisms which make them mice react to the smoke. The probabilistic regularity comes from the transfactual efficacy of the mouse organisms.
- Meteorology: the rainfall in a given month. It is very fortunate for the development of life on our planet that we have the chaotic alternation between cloud cover and clear sky, instead of a continuous cloud cover as in Venus or a continuous clear sky. Butterfly effect all over again, but it is possible to make probabilistic predictions since the fundamentals remain stable: the transfactual efficacy of the energy received from the sun and radiated back out into space.
- Markets: the total personal income in New York State in a given month. Market economies are a very much like the weather; planned economies would be more like production or life.
- Uncertainty: the exact date of Noah's birth. This is epistemic uncertainty: assuming that Noah was a real person, the date exists and we know a time range in which it must have been, but we do not know the details. Probabilistic methods can be used to represent this kind of uncertain knowledge, but other methods to represent this knowledge may be more appropriate.

- Indeterminacy: The closing of the Dow Jones Industrial Average (DJIA) or the temperature in New York City at 4 pm. on February 28, 2014: This is ontological uncertainty, not only epistemological uncertainty. Not only do we not know it, but it is objectively not yet decided what these data will be. Probability theory has limited applicability for the DJIA since it cannot be expected that the mechanisms determining the DJIA will be the same at that time, therefore we cannot base ourselves on the transfactual efficacy of some stable mechanisms. It is not known which stocks will be included in the DJIA at that time, or whether the US dollar will still be the world reserve currency and the New York stock exchange the pinnacle of international capital markets. Perhaps a different stock market index located somewhere else will at that time play the role the DJIA is playing today. We would not even be able to ask questions about that alternative index today.

Regarding the temperature, it is more defensible to assign a probability since the weather mechanisms have probably stayed the same, except for changes in global warming (unless mankind has learned by that time to manipulate the weather locally by cloud seeding etc.).

- Chaotic determinacy: the relative frequency of the digit 3 in the decimal representation of  $\pi$ : The laws by which the number  $\pi$  is defined have very little to do with the procedure by which numbers are expanded as decimals; therefore the former has no systematic influence on the latter. (It has a systematic influence, but not a systematic one; it is the error of actualism to think that every influence must be systematic.) But it is also known that laws of nature have remote effects: one of the most amazing theorems in mathematics is the formula  $\frac{\pi}{4} = 1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \dots$  which establishes a connection between the geometry of the circle and some simple arithmetics.
- Quantum mechanics: the proportion of photons absorbed by a polarizing filter: If these photons are already polarized (but in a different direction than the filter) then this is not epistemic uncertainty but ontological indeterminacy, since the polarized photons form a pure state, which is atomic in the algebra of events. In this case, the distinction between epistemic uncertainty and ontological indeterminacy is operational: the two alternatives follow different mathematics.
- Statistical mechanics: the velocity distribution of molecules in a gas at a given pressure and temperature. Thermodynamics cannot be reduced to the mechanics of molecules, since mechanics is reversible in time, while thermodynamics is not. An additional element is needed, which can be modeled using probability.



PROBLEM 2. *Not every kind of uncertainty can be formulated stochastically. Which other methods are available if stochastic means are inappropriate?*

ANSWER. Dialectics.  $\square$

PROBLEM 3. *How are the probabilities of rain in weather forecasts to be interpreted?*

ANSWER. Renyi in [Rén70, pp. 33/4]: “By saying that the probability of rain tomorrow is 80% (or, what amounts to the same, 0.8) the meteorologist means that in a situation similar to that observed on the given day, there is usually rain on the next day in about 8 out of 10 cases; thus, while it is not certain that it will rain tomorrow, the *degree of certainty* of this event is 0.8.”  $\square$

Pure uncertainty is as hard to generate as pure certainty; it is needed for encryption and numerical methods.

Here is an encryption scheme which leads to a random looking sequence of numbers (see [Rao97, p. 13]): First a string of binary random digits is generated which is known only to the sender and receiver. The sender converts his message into a string of binary digits. He then places the message string below the key string and obtains a coded string by changing every message bit to its alternative at all places where the key bit is 1 and leaving the others unchanged. The coded string which appears to be a random binary sequence is transmitted. The received message is decoded by making the changes in the same way as in encrypting using the key string which is known to the receiver.

PROBLEM 4. *Why is it important in the above encryption scheme that the key string is purely random and does not have any regularities?*

PROBLEM 5. [Knu81, pp. 7, 452] *Suppose you wish to obtain a decimal digit at random, not using a computer. Which of the following methods would be suitable?*

• a. *Open a telephone directory to a random place (i.e., stick your finger in it somewhere) and use the unit digit of the first number found on the selected page.*

ANSWER. This will often fail, since users select “round” numbers if possible. In some areas, telephone numbers are perhaps assigned randomly. But it is a mistake in any case to try to get several successive random numbers from the same page, since many telephone numbers are listed several times in a sequence.  $\square$

• b. *Same as a, but use the units digit of the page number.*

ANSWER. But do you use the left-hand page or the right-hand page? Say, use the left-hand page, divide by 2, and use the units digit.  $\square$

• c. *Roll a die which is in the shape of a regular icosahedron, whose twenty faces have been labeled with the digits 0, 0, 1, 1, . . . , 9, 9. Use the digit which appears on top, when the die comes to rest. (A felt table with a hard surface is recommended for rolling dice.)*

ANSWER. The markings on the face will slightly bias the die, but for practical purposes this method is quite satisfactory. See Math. Comp. 15 (1961), 94–95, for further discussion of this method.

• d. *Expose a geiger counter to a source of radioactivity for one minute (shield yourself) and use the unit digit of the resulting count. (Assume that the geiger counter displays the number of counts in decimal notation, and that the count is initially zero.)*

ANSWER. This is a difficult question thrown in purposely as a surprise. The number is uniformly distributed! One sees this best if one imagines the source of radioactivity is very low level, so that only a few emissions can be expected during this minute. If the average number of emissions per minute is  $\lambda$ , the probability that the counter registers  $k$  is  $e^{-\lambda}\lambda^k/k!$  (the Poisson distribution). So the digit 0 is selected with probability  $e^{-\lambda}\sum_{k=0}^{\infty}\lambda^{10k}/(10k)!$ , etc.

• e. *Glance at your wristwatch, and if the position of the second-hand is between  $6n$  and  $6(n+1)$ , choose the digit  $n$ .*

ANSWER. Okay, provided that the time since the last digit selected in this way is random. A bias may arise if borderline cases are not treated carefully. A better device seems to be to use a stopwatch which has been started long ago, and which one stops arbitrarily, and then one has the time necessary to read the display.

• f. *Ask a friend to think of a random digit, and use the digit he names.*

ANSWER. No, people usually think of certain digits (like 7) with higher probability.

• g. *Assume 10 horses are entered in a race and you know nothing whatever about their qualifications. Assign to these horses the digits 0 to 9, in arbitrary fashion, and after the race use the winner’s digit.*

ANSWER. Okay; your assignment of numbers to the horses had probability 1/10 of assigning a given digit to a winning horse.

## 2.2. Events as Sets

With every situation with uncertain outcome we associate its *sample space*  $U$  which represents the set of all possible outcomes (described by the characteristics which we are interested in).

*Events* are associated with subsets of the sample space, i.e., with bundles of outcomes that are observable in the given experimental setup. The set of all events we denote with  $\mathcal{F}$ . ( $\mathcal{F}$  is a set of subsets of  $U$ .)

Look at the example of rolling a die.  $U = \{1, 2, 3, 4, 5, 6\}$ . The events of getting an even number is associated with the subset  $\{2, 4, 6\}$ ; getting a six with  $\{6\}$ ; not getting a six with  $\{1, 2, 3, 4, 5\}$ , etc. Now look at the example of rolling two indistinguishable dice. Observable events may be: getting two ones, getting a one and a two, etc. But we cannot distinguish between the first die getting a one and the second

two, and vice versa. I.e., if we define the sample set to be  $U = \{1, \dots, 6\} \times \{1, \dots, 6\}$ , i.e., the set of all pairs of numbers between 1 and 6, then certain subsets are not observable.  $\{(1, 5)\}$  is not observable (unless the dice are marked or have different colors etc.), only  $\{(1, 5), (5, 1)\}$  is observable.

If the experiment is measuring the height of a person in meters, and we make the idealized assumption that the measuring instrument is infinitely accurate, then all possible outcomes are numbers between 0 and 3, say. Sets of outcomes one is usually interested in are whether the height falls within a given interval; therefore all intervals within the given range represent observable events.

If the sample space is finite or countably infinite, very often *all* subsets are observable events. If the sample set contains an uncountable continuum, it is not desirable to consider all subsets as observable events. Mathematically one can define quite crazy subsets which have no practical significance and which cannot be meaningfully given probabilities. For the purposes of Econ 7800, it is enough to say that all the subsets which we may reasonably define are candidates for observable events.

The “set of all possible outcomes” is well defined in the case of rolling a die and other games; but in social sciences, situations arise in which the outcome is open and the range of possible outcomes cannot be known beforehand. If one uses a probability theory based on the concept of a “set of possible outcomes” in such a situation, one reduces a process which is open and evolutionary to an imaginary predetermined and static “set.” Furthermore, in social theory, the mechanism by which these uncertain outcomes are generated are often internal to the members of the statistical population. The mathematical framework models these mechanisms as an extraneous “picking an element out of a pre-existing set.”

From given observable events we can derive new observable events by *set theoretical operations*. (All the operations below involve subsets of the same  $U$ .)

*Mathematical Note:* Notation of sets: there are two ways to denote a set: either by giving a rule, or by listing the elements. (The order in which the elements are listed, or the fact whether some elements are listed twice or not, is irrelevant.)

Here are the formal definitions of set theoretic operations. The letters  $A, B$ , etc. denote subsets of a given set  $U$  (events), and  $I$  is an arbitrary index set.  $\omega$  stands

for an element, and  $\omega \in A$  means that  $\omega$  is an element of  $A$ .

$$(2.2.1) \quad A \subset B \iff (\omega \in A \Rightarrow \omega \in B) \quad (A \text{ is contained in } B)$$

$$(2.2.2) \quad A \cap B = \{\omega : \omega \in A \text{ and } \omega \in B\} \quad (\text{intersection of } A \text{ and } B)$$

$$(2.2.3) \quad \bigcap_{i \in I} A_i = \{\omega : \omega \in A_i \text{ for all } i \in I\}$$

$$(2.2.4) \quad A \cup B = \{\omega : \omega \in A \text{ or } \omega \in B\} \quad (\text{union of } A \text{ and } B)$$

$$(2.2.5) \quad \bigcup_{i \in I} A_i = \{\omega : \text{there exists an } i \in I \text{ such that } \omega \in A_i\}$$

$$(2.2.6) \quad U \quad \text{Universal set: all } \omega \text{ we talk about are } \in U.$$

$$(2.2.7) \quad A' = \{\omega : \omega \notin A \text{ but } \omega \in U\}$$

$$(2.2.8) \quad \emptyset = \text{the empty set: } \omega \notin \emptyset \text{ for all } \omega.$$

These definitions can also be visualized by Venn diagrams; and for the purposes of this class, demonstrations with the help of Venn diagrams will be admissible in lieu of mathematical proofs.

**PROBLEM 6.** *For the following set-theoretical exercises it is sufficient that you draw the corresponding Venn diagrams and convince yourself by just looking at them that the statement is true. For those who are interested in a precise mathematical proof derived from the definitions of  $A \cup B$  etc. given above, should remember that a proof of the set-theoretical identity  $A = B$  usually has the form: first you show that  $\omega \in A$  implies  $\omega \in B$ , and then you show the converse.*

- a. *Prove that  $A \cup B = B \iff A \cap B = A$ .*

**ANSWER.** If one draws the Venn diagrams, one can see that either side is true if and only if  $A \subset B$ . If one wants a more precise proof, the following proof by contradiction seems most illuminating: Assume the lefthand side does not hold, i.e., there exists a  $\omega \in A$  but  $\omega \notin B$ . Then  $\omega \notin A \cap B$ , i.e.,  $A \cap B \neq A$ . Now assume the righthand side does not hold, i.e., there is a  $\omega \in A \cap B$  with  $\omega \notin B$ . This  $\omega$  lies in  $A \cup B$  but not in  $B$ , i.e., the lefthand side does not hold either.

- b. *Prove that  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$*

**ANSWER.** If  $\omega \in A$  then it is clearly always in the righthand side and in the lefthand side. If there is therefore any difference between the righthand and the lefthand side, it must be for  $\omega \notin A$ : If  $\omega \notin A$  and it is still in the lefthand side then it must be in  $B \cap C$ , therefore it is also in the righthand side. If  $\omega \notin A$  and it is in the righthand side, then it must be both in  $B$  and in  $C$ , therefore it is in the lefthand side.

- c. *Prove that  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ .*

ANSWER. If  $\omega \notin A$  then it is clearly neither in the righthand side nor in the lefthand side. If there is therefore any difference between the righthand and the lefthand side, it must be for the  $\omega \in A$ : If  $\omega \in A$  and it is in the lefthand side then it must be in  $B \cup C$ , i.e., in  $B$  or in  $C$  or in both, therefore it is also in the righthand side. If  $\omega \in A$  and it is in the righthand side, then it must be in either  $B$  or  $C$  or both, therefore it is in the lefthand side.  $\square$

- d. Prove that  $A \cap \left(\bigcup_{i=1}^{\infty} B_i\right) = \bigcup_{i=1}^{\infty} (A \cap B_i)$ .

ANSWER. Proof: If  $\omega$  in lefthand side, then it is in  $A$  and in at least one of the  $B_i$ , say it is in  $B_k$ . Therefore it is in  $A \cap B_k$ , and therefore it is in the righthand side. Now assume, conversely, that  $\omega$  is in the righthand side; then it is at least in one of the  $A \cap B_i$ , say it is in  $A \cap B_k$ . Hence it is in  $A$  and in  $B_k$ , i.e., in  $A$  and in  $\bigcup B_i$ , i.e., it is in the lefthand side.  $\square$

PROBLEM 7. 3 points Draw a Venn Diagram which shows the validity of de Morgan's laws:  $(A \cup B)' = A' \cap B'$  and  $(A \cap B)' = A' \cup B'$ . If done right, the same Venn diagram can be used for both proofs.

ANSWER. There is a proof in [HT83, p. 12]. Draw  $A$  and  $B$  inside a box which represents  $U$ , and shade  $A'$  from the left (blue) and  $B'$  from the right (yellow), so that  $A' \cap B'$  is cross shaded (green); then one can see these laws.  $\square$

PROBLEM 8. 3 points [HT83, Exercise 1.2-13 on p. 14] Evaluate the following unions and intersections of intervals. Use the notation  $(a, b)$  for open and  $[a, b]$  for closed intervals,  $(a, b]$  or  $[a, b)$  for half open intervals,  $\{a\}$  for sets containing one element only, and  $\emptyset$  for the empty set.

$$(2.2.9) \quad \bigcup_{n=1}^{\infty} \left(\frac{1}{n}, 2\right) = \bigcap_{n=1}^{\infty} \left(0, \frac{1}{n}\right) =$$

$$(2.2.10) \quad \bigcup_{n=1}^{\infty} \left[\frac{1}{n}, 2\right] = \bigcap_{n=1}^{\infty} \left[0, 1 + \frac{1}{n}\right] =$$

ANSWER.

$$(2.2.11) \quad \bigcup_{n=1}^{\infty} \left(\frac{1}{n}, 2\right) = (0, 2) \quad \bigcap_{n=1}^{\infty} \left(0, \frac{1}{n}\right) = \emptyset$$

$$(2.2.12) \quad \bigcup_{n=1}^{\infty} \left[\frac{1}{n}, 2\right] = (0, 2] \quad \bigcap_{n=1}^{\infty} \left[0, 1 + \frac{1}{n}\right] = [0, 1]$$

Explanation of  $\bigcup_{n=1}^{\infty} \left[\frac{1}{n}, 2\right]$ : for every  $\alpha$  with  $0 < \alpha \leq 2$  there is a  $n$  with  $\frac{1}{n} \leq \alpha$ , but 0 itself is in none of the intervals.  $\square$

The set operations become logical operations if applied to events. Every experiment returns an element  $\omega \in U$  as outcome. Here  $\omega$  is rendered green in the electronic version of these notes (and in an upright font in the version for black-and-white printouts), because  $\omega$  does not denote a specific element of  $U$ , but it depends

on chance which element is picked. I.e., the green color (or the unusual font) indicates that  $\omega$  is “alive.” We will also render the events themselves (as opposed to their set-theoretical counterparts) in green (or in an upright font).

- We say that the event  $A$  has *occurred* when  $\omega \in A$ .
- If  $A \subset B$  then event  $A$  implies event  $B$ , and we will write this directly in terms of events as  $A \subset B$ .
- The set  $A \cap B$  is associated with the event that both  $A$  and  $B$  occur (e.g., an even number smaller than six), and considered as an event, not a set. The event that both  $A$  and  $B$  occur will be written  $A \cap B$ .
- Likewise,  $A \cup B$  is the event that either  $A$  or  $B$ , or both, occur.
- $A'$  is the event that  $A$  does not occur.
- $\mathbf{U}$  the event that always occurs (as long as one performs the experiment).
- The empty set  $\emptyset$  is associated with the impossible event  $\emptyset$ , because whatever the value  $\omega$  of the chance outcome  $\omega$  of the experiment, it is always  $\omega \notin \emptyset$ .

If  $A \cap B = \emptyset$ , the set theoretician calls  $A$  and  $B$  “disjoint,” and the probabilistic theoretician calls the events  $A$  and  $B$  “mutually exclusive.” If  $A \cup B = \mathbf{U}$ , then  $A$  and  $B$  are called “collectively exhaustive.”

The set  $\mathcal{F}$  of all observable events must be a  $\sigma$ -algebra, i.e., it must satisfy:

$$\emptyset \in \mathcal{F}$$

$$A \in \mathcal{F} \Rightarrow A' \in \mathcal{F}$$

$$A_1, A_2, \dots \in \mathcal{F} \Rightarrow A_1 \cup A_2 \cup \dots \in \mathcal{F} \quad \text{which can also be written as } \bigcup_{i=1,2,\dots} A_i \in \mathcal{F}$$

$$A_1, A_2, \dots \in \mathcal{F} \Rightarrow A_1 \cap A_2 \cap \dots \in \mathcal{F} \quad \text{which can also be written as } \bigcap_{i=1,2,\dots} A_i \in \mathcal{F}$$

### 2.3. The Axioms of Probability

A probability measure  $\Pr : \mathcal{F} \rightarrow \mathbb{R}$  is a mapping which assigns to every event a number, the probability of this event. This assignment must be compatible with the set-theoretic operations between events in the following way:

$$(2.3.1) \quad \Pr[\mathbf{U}] = 1$$

$$(2.3.2) \quad \Pr[A] \geq 0 \quad \text{for all events } A$$

$$(2.3.3) \quad \text{If } A_i \cap A_j = \emptyset \text{ for all } i, j \text{ with } i \neq j \text{ then } \Pr\left[\bigcup_{i=1}^{\infty} A_i\right] = \sum_{i=1}^{\infty} \Pr[A_i]$$

Here an infinite sum is mathematically defined as the limit of partial sums. The axioms make probability what mathematicians call a *measure*, like area or weight.

In a Venn diagram, one might therefore interpret the probability of the events as the *area* of the bubble representing the event.

PROBLEM 9. *Prove that  $\Pr[A'] = 1 - \Pr[A]$ .*

ANSWER. Follows from the fact that  $A$  and  $A'$  are disjoint and their union  $U$  has probability 1.  $\square$

PROBLEM 10. *2 points Prove that  $\Pr[A \cup B] = \Pr[A] + \Pr[B] - \Pr[A \cap B]$ .*

ANSWER. For Econ 7800 it is sufficient to argue it out intuitively: if one adds  $\Pr[A] + \Pr[B]$  then one counts  $\Pr[A \cap B]$  twice and therefore has to subtract it again.

The brute force mathematical proof guided by this intuition is somewhat verbose: Define  $D = A \cap B'$ ,  $E = A \cap B$ , and  $F = A' \cap B$ .  $D$ ,  $E$ , and  $F$  satisfy

$$(2.3.4) \quad D \cup E = (A \cap B') \cup (A \cap B) = A \cap (B' \cup B) = A \cap U = A,$$

$$(2.3.5) \quad E \cup F = B,$$

$$(2.3.6) \quad D \cup E \cup F = A \cup B.$$

You may need some of the properties of unions and intersections in Problem 6. Next step is to prove that  $D$ ,  $E$ , and  $F$  are mutually exclusive. Therefore it is easy to take probabilities

$$(2.3.7) \quad \Pr[A] = \Pr[D] + \Pr[E];$$

$$(2.3.8) \quad \Pr[B] = \Pr[E] + \Pr[F];$$

$$(2.3.9) \quad \Pr[A \cup B] = \Pr[D] + \Pr[E] + \Pr[F].$$

Take the sum of (2.3.7) and (2.3.8), and subtract (2.3.9):

$$(2.3.10) \quad \Pr[A] + \Pr[B] - \Pr[A \cup B] = \Pr[E] = \Pr[A \cap B];$$

A shorter but trickier alternative proof is the following. First note that  $A \cup B = A \cup (A' \cap B)$  and that this is a disjoint union, i.e.,  $\Pr[A \cup B] = \Pr[A] + \Pr[A' \cap B]$ . Then note that  $B = (A \cap B) \cup (A' \cap B)$ , and this is a disjoint union, therefore  $\Pr[B] = \Pr[A \cap B] + \Pr[A' \cap B]$ , or  $\Pr[A' \cap B] = \Pr[B] - \Pr[A \cap B]$ . Putting this together gives the result.  $\square$

PROBLEM 11. *1 point Show that for arbitrary events  $A$  and  $B$ ,  $\Pr[A \cup B] \leq \Pr[A] + \Pr[B]$ .*

ANSWER. From Problem 10 we know that  $\Pr[A \cup B] = \Pr[A] + \Pr[B] - \Pr[A \cap B]$ , and from axiom (2.3.2) follows  $\Pr[A \cap B] \geq 0$ .  $\square$

PROBLEM 12. *2 points (Bonferroni inequality) Let  $A$  and  $B$  be two events. Write  $\Pr[A] = 1 - \alpha$  and  $\Pr[B] = 1 - \beta$ , show that  $\Pr[A \cap B] \geq 1 - (\alpha + \beta)$ . You are allowed to use that  $\Pr[A \cup B] = \Pr[A] + \Pr[B] - \Pr[A \cap B]$  (Problem 10), and that all probabilities are  $\leq 1$ .*

ANSWER.

$$(2.3.11) \quad \Pr[A \cup B] = \Pr[A] + \Pr[B] - \Pr[A \cap B] \leq 1$$

$$(2.3.12) \quad \Pr[A] + \Pr[B] \leq 1 + \Pr[A \cap B]$$

$$(2.3.13) \quad \Pr[A] + \Pr[B] - 1 \leq \Pr[A \cap B]$$

$$(2.3.14) \quad 1 - \alpha + 1 - \beta - 1 = 1 - \alpha - \beta \leq \Pr[A \cap B]$$

PROBLEM 13. *(Not eligible for in-class exams) Given a rising sequence of events  $B_1 \subset B_2 \subset B_3 \cdots$ , define  $B = \bigcup_{i=1}^{\infty} B_i$ . Show that  $\Pr[B] = \lim_{i \rightarrow \infty} \Pr[B_i]$ .*

ANSWER. Define  $C_1 = B_1$ ,  $C_2 = B_2 \cap B'_1$ ,  $C_3 = B_3 \cap B'_2$ , etc. Then  $C_i \cap C_j = \emptyset$  for  $i \neq j$ , and  $B_n = \bigcup_{i=1}^n C_i$  and  $B = \bigcup_{i=1}^{\infty} C_i$ . In other words, now we have represented every  $B_n$  and  $B$  as a union of disjoint sets, and can therefore apply the third probability axiom (2.3.3):  $\Pr[B] = \sum_{i=1}^{\infty} \Pr[C_i]$ . The infinite sum is merely a short way of writing  $\Pr[B] = \lim_{n \rightarrow \infty} \sum_{i=1}^n \Pr[C_i]$ , i.e., the infinite sum is the limit of the finite sums. But since these finite sums are exactly  $\sum_{i=1}^n \Pr[C_i] = \Pr[\bigcup_{i=1}^n C_i] = \Pr[B_n]$ , the assertion follows. This proof, as it stands, is for our purposes entirely acceptable. One can make some steps in this proof still more stringent. For instance, one might use induction to prove  $B_n = \bigcup_{i=1}^n C_i$ . And how does one show that  $B = \bigcup_{i=1}^{\infty} C_i$ ? Well, one knows that  $C_i \subset B$  therefore  $\bigcup_{i=1}^{\infty} C_i \subset \bigcup_{i=1}^{\infty} B_i = B$ . Now take an  $\omega \in B$ . Then it lies in at least one of the  $B_i$ , it can be in many of them. Let  $k$  be the smallest  $k$  for which  $\omega \in B_k$ . If  $k = 1$ , then  $\omega \in C_1 = B_1$  as well. Otherwise,  $\omega \notin B_{k-1}$ , and therefore  $\omega \in C_k$ . I.e., any element in  $B$  lies in at least one of the  $C_k$ , therefore  $B \subset \bigcup_{i=1}^{\infty} C_i$ .

PROBLEM 14. *(Not eligible for in-class exams) From problem 13 derive a theorem: the following: if  $A_1 \supset A_2 \supset A_3 \cdots$  is a declining sequence, and  $A = \bigcap_i A_i$ , then  $\Pr[A] = \lim \Pr[A_i]$ .*

ANSWER. If the  $A_i$  are declining, then their complements  $B_i = A'_i$  are rising:  $B_1 \subset B_2 \subset B_3 \cdots$  are rising; therefore I know the probability of  $B = \bigcup B_i$ . Since by de Morgan's laws,  $B = A'$ , this gives me also the probability of  $A$ .

The results regarding the probabilities of rising or declining sequences are equivalent to the third probability axiom. This third axiom can therefore be considered as a continuity condition for probabilities.

If  $U$  is finite or countably infinite, then the probability measure is uniquely determined if one knows the probability of every one-element set. We will call  $\Pr[\{\omega\}] = p(\omega)$  the probability mass function. Other terms used for it in the literature are probability function, or even probability density function (although it is not a density, more about this below). If  $U$  has more than countably infinite elements, the probabilities of one-element sets may not give enough information to define the whole probability measure.

*Mathematical Note:* Not all infinite sets are countable. Here is a proof, by contradiction, that the real numbers between 0 and 1 are not countable: assume there is an enumeration, i.e., a sequence  $a_1, a_2, \dots$  which contains them all. We

them underneath each other in their (possibly infinite) decimal representation, where  $0.d_{i1}d_{i2}d_{i3}\dots$  is the decimal representation of  $a_i$ . Then any real number whose decimal representation is such that the first digit is *not* equal to  $d_{11}$ , the second digit is *not* equal  $d_{22}$ , the third *not* equal  $d_{33}$ , etc., is a real number which is *not* contained in this enumeration. That means, an enumeration which contains all real numbers cannot exist.

On the real numbers between 0 and 1, the length measure (which assigns to each interval its length, and to sets composed of several intervals the sums of the lengths, etc.) is a probability measure. In this probability field, every one-element subset of the sample set has zero probability.

This shows that events other than  $\emptyset$  may have zero probability. In other words, if an event has probability 0, this does not mean it is logically impossible. It may well happen, but it happens so infrequently that in repeated experiments the *average* number of occurrences converges toward zero.

## 2.4. Objective and Subjective Interpretation of Probability

The mathematical probability axioms apply to both objective and subjective interpretation of probability.

The *objective* interpretation considers probability a quasi physical property of the experiment. One cannot simply say:  $\Pr[A]$  is the relative frequency of the occurrence of  $A$ , because we know intuitively that this frequency does not necessarily converge. E.g., even with a fair coin it is physically possible that one always gets head, or that one gets some other sequence which does not converge towards  $\frac{1}{2}$ . The above axioms resolve this dilemma, because they allow to derive the theorem that the relative frequencies converges towards the probability *with probability one*.

*Subjectivist* interpretation (de Finetti: “probability does not exist”) defines probability in terms of people’s ignorance and willingness to take bets. Interesting for economists because it uses money and utility, as in expected utility. Call “a lottery on  $A$ ” a lottery which pays \$1 if  $A$  occurs, and which pays nothing if  $A$  does not occur. If a person is willing to pay  $p$  dollars for a lottery on  $A$  and  $1 - p$  dollars for a lottery on  $A'$ , then, according to a subjectivist definition of probability, he assigns subjective probability  $p$  to  $A$ .

There is the presumption that his willingness to bet does not depend on the size of the payoff (i.e., the payoffs are considered to be small amounts).

**PROBLEM 15.** Assume  $A$ ,  $B$ , and  $C$  are a complete disjunction of events, i.e., they are mutually exclusive and  $A \cup B \cup C = U$ , the universal set.

• a. 1 point Arnold assigns subjective probability  $p$  to  $A$ ,  $q$  to  $B$ , and  $r$  to  $C$ . Explain exactly what this means.

**ANSWER.** We know six different bets which Arnold is always willing to make, not only on  $B$ , and  $C$ , but also on their complements.

• b. 1 point Assume that  $p + q + r > 1$ . Name three lotteries which Arnold would be willing to buy, the net effect of which would be that he loses with certainty.

**ANSWER.** Among those six we have to pick subsets that make him a sure loser. If  $p + q + r > 1$  then we sell him a bet on  $A$ , one on  $B$ , and one on  $C$ . The payoff is always 1, and the cost is  $p + q + r > 1$ .

• c. 1 point Now assume that  $p + q + r < 1$ . Name three lotteries which Arnold would be willing to buy, the net effect of which would be that he loses with certainty.

**ANSWER.** If  $p + q + r < 1$ , then we sell him a bet on  $A'$ , one on  $B'$ , and one on  $C'$ . The payoff is 2, and the cost is  $1 - p + 1 - q + 1 - r > 2$ .

• d. 1 point Arnold is therefore only coherent if  $\Pr[A] + \Pr[B] + \Pr[C] = 1$ . Show that the additivity of probability can be derived from coherence, i.e., show that a subjective probability that satisfies the rule: whenever  $A$ ,  $B$ , and  $C$  is a complete disjunction of events, then the sum of their probabilities is 1, is additive, i.e.,  $\Pr[A \cup B] = \Pr[A] + \Pr[B]$ .

**ANSWER.** Since  $r$  is his subjective probability of  $C$ ,  $1 - r$  must be his subjective probability of  $C' = A \cup B$ . Since  $p + q + r = 1$ , it follows  $1 - r = p + q$ .

This last problem indicates that the finite additivity axiom follows from the requirement that the bets be consistent or, as subjectivists say, “coherent” with each other. However, it is not possible to derive the additivity for countably infinite sequences of events from such an argument.

## 2.5. Counting Rules

In this section we will be working in a *finite* probability space, in which all atomic events have equal probabilities. The acts of rolling dice or drawing balls from urns can be modeled by such spaces. In order to compute the probability of a given event one must *count* the elements of the set which this event represents. In other words we *count* how many different ways there are to achieve a certain outcome. This can be tricky, and we will develop some general principles how to do it.

**PROBLEM 16.** You throw two dice.

• a. 1 point What is the probability that the sum of the numbers shown is five or less?

**ANSWER.**  $\frac{11}{31} \frac{12}{32} \frac{13}{33} \frac{14}{34}$ , i.e., 10 out of 36 possibilities, gives the probability  $\frac{5}{18}$ .

• b. 1 point What is the probability that both of the numbers shown are five or less?



The binomial coefficients also occur in the Binomial Theorem

$$(2.6.3) \quad (a + b)^n = a^n + \binom{n}{1}a^{n-1}b + \dots + \binom{n}{n-1}ab^{n-1} + b^n = \sum_{k=0}^n \binom{n}{k}a^{n-k}b^k$$

Why? When the  $n$  factors  $a + b$  are multiplied out, each of the resulting terms selects from each of the  $n$  original factors either  $a$  or  $b$ . The term  $a^{n-k}b^k$  occurs therefore  $\binom{n}{n-k} = \binom{n}{k}$  times.

As an application: If you set  $a = 1$ ,  $b = 1$ , you simply get a sum of binomial coefficients, i.e., you get the number of subsets in a set with  $n$  elements: it is  $2^n$  (always count the empty set as one of the subsets). The number of all subsets is easily counted directly. You go through the set element by element and about every element you ask: is it in the subset or not? I.e., for every element you have two possibilities, therefore by the multiplication principle the total number of possibilities is  $2^n$ .

### 2.7. Conditional Probability

The concept of conditional probability is arguably more fundamental than probability itself. Every probability is conditional, since we must know that the “experiment” has happened before we can speak of probabilities. [Ame94, p. 10] and [Rén70] give axioms for conditional probability which take the place of the above axioms (2.3.1), (2.3.2) and (2.3.3). However we will follow here the common procedure of defining conditional probabilities in terms of the unconditional probabilities:

$$(2.7.1) \quad \Pr[B|A] = \frac{\Pr[B \cap A]}{\Pr[A]}$$

How can we motivate (2.7.1)? If we know that  $A$  has occurred, then of course the only way that  $B$  occurs is when  $B \cap A$  occurs. But we want to multiply all probabilities of subsets of  $A$  with an appropriate proportionality factor so that the probability of the event  $A$  itself becomes  $= 1$ .

**PROBLEM 20.** 3 points Let  $A$  be an event with nonzero probability. Show that the probability conditionally on  $A$ , i.e., the mapping  $B \mapsto \Pr[B|A]$ , satisfies all the axioms of a probability measure:

$$(2.7.2) \quad \Pr[U|A] = 1$$

$$(2.7.3) \quad \Pr[B|A] \geq 0 \quad \text{for all events } B$$

$$(2.7.4) \quad \Pr\left[\bigcup_{i=1}^{\infty} B_i|A\right] = \sum_{i=1}^{\infty} \Pr[B_i|A] \quad \text{if } B_i \cap B_j = \emptyset \text{ for all } i, j \text{ with } i \neq j.$$

**ANSWER.**  $\Pr[U|A] = \Pr[U \cap A]/\Pr[A] = 1$ .  $\Pr[B|A] = \Pr[B \cap A]/\Pr[A] \geq 0$  because  $\Pr[B \cap A]$  and  $\Pr[A] > 0$ . Finally,

$$(2.7.5) \quad \Pr\left[\bigcup_{i=1}^{\infty} B_i|A\right] = \frac{\Pr[(\bigcup_{i=1}^{\infty} B_i) \cap A]}{\Pr[A]} = \frac{\Pr[\bigcup_{i=1}^{\infty} (B_i \cap A)]}{\Pr[A]} = \frac{1}{\Pr[A]} \sum_{i=1}^{\infty} \Pr[B_i \cap A] = \sum_{i=1}^{\infty} \Pr[B_i|A]$$

First equal sign is definition of conditional probability, second is distributivity of unions and intersections (Problem 6 d), third because the  $B_i$  are disjoint and therefore the  $B_i \cap A$  are even more disjoint:  $B_i \cap A \cap B_j \cap A = B_i \cap B_j \cap A = \emptyset \cap A = \emptyset$  for all  $i, j$  with  $i \neq j$ , and the last equal sign again by the definition of conditional probability.

**PROBLEM 21.** You draw two balls without replacement from an urn which has 10 white and 14 black balls.

If both balls are white, you roll a die, and your payoff is the number which the die shows in dollars.

If one ball is black and one is white, you flip a coin until you get your first head, and your payoff will be the number of flips it takes you to get a head, in dollars again.

If both balls are black, you draw from a deck of 52 cards, and you get the number shown on the card in dollars. (Ace counts as one, J, Q, and K as 11, 12, 13, i.e., basically the deck contains every number between 1 and 13 four times.)

Show that the probability that you receive exactly two dollars in this game is  $\frac{1}{6}$ .

**ANSWER.** You know a complete disjunction of events:  $U = \{ww\} \cup \{bb\} \cup \{wb\}$ , with  $\Pr[\{ww\}] = \frac{7}{21} \cdot \frac{6}{20} = \frac{1}{10}$ ;  $\Pr[\{bb\}] = \frac{14}{21} \cdot \frac{13}{20} = \frac{13}{30}$ ;  $\Pr[\{bw\}] = \frac{7}{21} \cdot \frac{14}{20} + \frac{14}{21} \cdot \frac{7}{20} = \frac{7}{15}$ . Furthermore you know the conditional probabilities of getting 2 dollars conditionally on each of these events:  $\Pr[\{2\}|\{ww\}] = \frac{1}{6}$ ;  $\Pr[\{2\}|\{bb\}] = \frac{1}{13}$ ;  $\Pr[\{2\}|\{wb\}] = \frac{1}{4}$ . Now  $\Pr[\{2\} \cap \{ww\}] = \Pr[\{2\}|\{ww\}] \Pr[\{ww\}]$  etc., therefore

$$(2.7.6) \quad \Pr[\{2\}] = \Pr[\{2\} \cap \{ww\}] + \Pr[\{2\} \cap \{bw\}] + \Pr[\{2\} \cap \{bb\}]$$

$$(2.7.7) \quad = \frac{1}{6} \frac{7}{21} \frac{6}{20} + \frac{1}{4} \left( \frac{7}{21} \frac{14}{20} + \frac{14}{21} \frac{7}{20} \right) + \frac{1}{13} \frac{14}{21} \frac{13}{20}$$

$$(2.7.8) \quad = \frac{1}{6} \frac{1}{10} + \frac{1}{4} \frac{7}{15} + \frac{1}{13} \frac{13}{30} = \frac{1}{6}$$

**PROBLEM 22.** 2 points  $A$  and  $B$  are arbitrary events. Prove that the probability of  $B$  can be written as:

$$(2.7.9) \quad \Pr[B] = \Pr[B|A] \Pr[A] + \Pr[B|A'] \Pr[A']$$

This is the law of iterated expectations (6.6.2) in the case of discrete random variables: it might be written as  $\Pr[B] = E[\Pr[B|A]]$ .

**ANSWER.**  $B = B \cap U = B \cap (A \cup A') = (B \cap A) \cup (B \cap A')$  and this union is disjoint, so  $\Pr[B] = \Pr[B \cap A] + \Pr[B \cap A']$ . Therefore  $\Pr[B] = \Pr[B \cap A] + \Pr[B \cap A']$ . Now apply definition of conditional probability to get  $\Pr[B \cap A] = \Pr[B|A] \Pr[A]$  and  $\Pr[B \cap A'] = \Pr[B|A'] \Pr[A']$ .

PROBLEM 23. 2 points Prove the following lemma: If  $\Pr[B|A_1] = \Pr[B|A_2]$  (call it  $c$ ) and  $A_1 \cap A_2 = \emptyset$  (i.e.,  $A_1$  and  $A_2$  are disjoint), then also  $\Pr[B|A_1 \cup A_2] = c$ .

ANSWER.

$$\begin{aligned} \Pr[B|A_1 \cup A_2] &= \frac{\Pr[B \cap (A_1 \cup A_2)]}{\Pr[A_1 \cup A_2]} = \frac{\Pr[(B \cap A_1) \cup (B \cap A_2)]}{\Pr[A_1 \cup A_2]} \\ (2.7.10) \quad &= \frac{\Pr[B \cap A_1] + \Pr[B \cap A_2]}{\Pr[A_1] + \Pr[A_2]} = \frac{c\Pr[A_1] + c\Pr[A_2]}{\Pr[A_1] + \Pr[A_2]} = c. \end{aligned}$$

□

PROBLEM 24. Show by counterexample that the requirement  $A_1 \cap A_2 = \emptyset$  is necessary for this result to hold. Hint: use the example in Problem 38 with  $A_1 = \{HH, HT\}$ ,  $A_2 = \{HH, TH\}$ ,  $B = \{HH, TT\}$ .

ANSWER.  $\Pr[B|A_1] = 1/2$  and  $\Pr[B|A_2] = 1/2$ , but  $\Pr[B|A_1 \cup A_2] = 1/3$ . □

The conditional probability can be used for computing probabilities of intersections of events.

PROBLEM 25. [Lar82, exercises 2.5.1 and 2.5.2 on p. 57, solutions on p. 597, but no discussion]. Five white and three red balls are laid out in a row at random.

• a. 3 points What is the probability that both end balls are white? What is the probability that one end ball is red and the other white?

ANSWER. You can lay the first ball first and the last ball second: for white balls, the probability is  $\frac{5}{8} \cdot \frac{4}{7} = \frac{5}{14}$ ; for one white, one red it is  $\frac{5}{8} \cdot \frac{3}{7} + \frac{3}{8} \cdot \frac{5}{7} = \frac{15}{28}$ . □

• b. 4 points What is the probability that all red balls are together? What is the probability that all white balls are together?

ANSWER. All red balls together is the same as 3 reds first, multiplied by 6, because you may have between 0 and 5 white balls before the first red.  $\frac{3}{8} \cdot \frac{2}{7} \cdot \frac{1}{6} \cdot 6 = \frac{3}{28}$ . For the white balls you get  $\frac{5}{8} \cdot \frac{4}{7} \cdot \frac{3}{6} \cdot \frac{2}{5} \cdot \frac{1}{4} \cdot 4 = \frac{1}{14}$ .

BTW, 3 reds first is same probability as 3 reds last, i.e., the 5 whites first:  $\frac{5}{8} \cdot \frac{4}{7} \cdot \frac{3}{6} \cdot \frac{2}{5} \cdot \frac{1}{4} = \frac{3}{8} \cdot \frac{2}{7} \cdot \frac{1}{6}$ . □

PROBLEM 26. The first three questions here are discussed in [Lar82, example 2.6.3 on p. 62]: There is an urn with 4 white and 8 black balls. You take two balls out without replacement.

• a. 1 point What is the probability that the first ball is white?

ANSWER.  $1/3$  □

• b. 1 point What is the probability that both balls are white?

ANSWER. It is  $\Pr[\text{second ball white}|\text{first ball white}]\Pr[\text{first ball white}] = \frac{3}{3+8} \cdot \frac{4}{4+8} = \frac{1}{11}$ . □

• c. 1 point What is the probability that the second ball is white?

ANSWER. It is  $\Pr[\text{first ball white and second ball white}] + \Pr[\text{first ball black and second ball white}]$

$$(2.7.11) \quad = \frac{3}{3+8} \cdot \frac{4}{4+8} + \frac{4}{7+4} \cdot \frac{8}{8+4} = \frac{1}{3}.$$

This is the same as the probability that the first ball is white. The probabilities are not dependent on the order in which one takes the balls out. This property is called “exchangeability.” One can also see it this way: Assume you number the balls at random, from 1 to 12. Then the probability for a white ball to have the number 2 assigned to it is obviously  $\frac{1}{3}$ .

• d. 1 point What is the probability that both of them are black?

ANSWER.  $\frac{8}{12} \cdot \frac{7}{11} = \frac{2}{3} \cdot \frac{7}{11} = \frac{14}{33}$  (or  $\frac{56}{132}$ ).

• e. 1 point What is the probability that both of them have the same color?

ANSWER. The sum of the two above,  $\frac{14}{33} + \frac{1}{11} = \frac{17}{33}$  (or  $\frac{68}{132}$ ).

Now you take three balls out without replacement.

• f. 2 points Compute the probability that at least two of the three balls are white.

ANSWER. It is  $\frac{13}{55}$ . The possibilities are  $wbw$ ,  $wbw$ ,  $bww$ , and  $www$ . Of the first three, each has probability  $\frac{4}{12} \cdot \frac{3}{11} \cdot \frac{8}{10}$ . Therefore the probability for exactly two being white is  $\frac{288}{1320} = \frac{12}{55}$ . The probability for  $www$  is  $\frac{4 \cdot 3 \cdot 2}{12 \cdot 11 \cdot 10} = \frac{24}{1320} = \frac{1}{55}$ . Add this to get  $\frac{312}{1320} = \frac{13}{55}$ . More systematically, answer is  $\left( \binom{4}{2} \binom{8}{1} + \binom{4}{1} \binom{8}{2} \right) / \binom{12}{3}$ .

• g. 1 point Compute the probability that at least two of the three are black.

ANSWER. It is  $\frac{42}{55}$ . For exactly two:  $\frac{672}{1320} = \frac{28}{55}$ . For three it is  $\frac{(8)(7)(6)}{(12)(11)(10)} = \frac{336}{1320} = \frac{14}{55}$ . Together  $\frac{1008}{1320} = \frac{42}{55}$ . One can also get it as: it is the complement of the last, or as  $\left( \binom{8}{2} \binom{4}{1} \right) / \binom{12}{3}$ .

• h. 1 point Compute the probability that two of the three are of the same color and the third of a different color.

ANSWER. It is  $\frac{960}{1320} = \frac{40}{55} = \frac{8}{11}$ , or  $\left( \binom{4}{1} \binom{8}{2} + \binom{4}{2} \binom{8}{1} \right) / \binom{12}{3}$ .

• i. 1 point Compute the probability that at least two of the three are of the same color.

ANSWER. This probability is 1. You have 5 black socks and 5 white socks in your drawer. There is a fire at night and you must get out of your apartment in two minutes. There is no light. You fumble in the dark for the drawer. How many socks do you have to take out so that you have at least 2 of the same color? The answer is 3 socks.

PROBLEM 27. If a poker hand of five cards is drawn from a deck, what is the probability that it will contain three aces? (How can the concept of conditional probability help in answering this question?)



ANSWER. [Ame94, example 2.3.3 on p. 9] and [Ame94, example 2.5.1 on p. 13] give two alternative ways to do it. The second answer uses conditional probability: Probability to draw three aces in a row first and then 2 nonaces is  $\frac{4}{52} \frac{3}{51} \frac{2}{50} \frac{48}{49} \frac{47}{48}$ . Then multiply this by  $\binom{5}{3} = \frac{5 \cdot 4 \cdot 3}{1 \cdot 2 \cdot 3} = 10$ . This gives 0.0017, i.e., 0.17%.  $\square$

PROBLEM 28. *2 points A friend tosses two coins. You ask: “did one of them land heads?” Your friend answers, “yes.” What’s the probability that the other also landed heads?*

ANSWER.  $U = \{HH, HT, TH, TT\}$ ; Probability is  $\frac{1}{4} / \frac{3}{4} = \frac{1}{3}$ .  $\square$

PROBLEM 29. *(Not eligible for in-class exams) [Ame94, p. 5] What is the probability that a person will win a game in tennis if the probability of his or her winning a point is  $p$ ?*

ANSWER.

$$(2.7.12) \quad p^4 \left( 1 + 4(1-p) + 10(1-p)^2 + \frac{20p(1-p)^3}{1-2p(1-p)} \right)$$

How to derive this:  $\{ssss\}$  has probability  $p^4$ ;  $\{ssfs\}$ ,  $\{sffs\}$ ,  $\{sfss\}$ , and  $\{fsss\}$  have probability  $4p^4(1-p)$ ;  $\{ssffs\}$  etc. (2  $f$  and 3  $s$  in the first 5, and then an  $s$ , together  $\binom{5}{2} = 10$  possibilities) have probability  $10p^4(1-p)^2$ . Now  $\{ssfff\}$  and  $\binom{6}{3} = 20$  other possibilities give deuce at least once in the game, i.e., the probability of deuce is  $20p^3(1-p)^3$ . Now  $\Pr[\text{win}|\text{deuce}] = p^2 + 2p(1-p)\Pr[\text{win}|\text{deuce}]$ , because you win either if you score twice in a row ( $p^2$ ) or if you get deuce again (probability  $2p(1-p)$ ) and then win. Solve this to get  $\Pr[\text{win}|\text{deuce}] = p^2 / (1 - 2p(1-p))$  and then multiply this conditional probability with the probability of getting deuce at least once:  $\Pr[\text{win after at least one deuce}] = 20p^3(1-p)^3 p^2 / (1 - 2p(1-p))$ . This gives the last term in (2.7.12).  $\square$

PROBLEM 30. *(Not eligible for in-class exams) Andy, Bob, and Chris play the following game: each of them draws a card without replacement from a deck of 52 cards. The one who has the highest card wins. If there is a tie (like: two kings and no aces), then that person wins among those who drew this highest card whose name comes first in the alphabet. What is the probability for Andy to be the winner? For Bob? For Chris? Does this probability depend on the order in which they draw their cards out of the stack?*

ANSWER. Let A be the event that Andy wins, B that Bob, and C that Chris wins.

One way to approach this problem is to ask: what are the chances for Andy to win when he draws a king?, etc., i.e., compute it for all 13 different cards. Then: what are the chances for Bob to win when he draws a king, and also his chances for the other cards, and then for Chris.

It is computationally easier to make the following partitioning of all outcomes: Either all three cards drawn are different (call this event D), or all three cards are equal (event E), or two of the three cards are equal (T). This third case will have to be split into  $T = H \cup L$ , according to whether the card that is different is higher or lower.

If all three cards are different, then Andy, Bob, and Chris have equal chances of winning; if all three cards are equal, then Andy wins. What about the case that two cards are the same and the

third is different? There are two possibilities. If the card that is different is higher than the two that are the same, then the chances of winning are evenly distributed; but if the two equal cards are higher, then Andy has a  $\frac{2}{3}$  chance of winning (when the distribution of the cards  $Y$  (low) and  $Z$  (higher) among  $ABC$  is  $ZZY$  and  $ZYZ$ ), and Bob has a  $\frac{1}{3}$  chance of winning (when the distribution is  $YZZ$ ). What we just did was computing the conditional probabilities  $\Pr[A|E]$ ,  $\Pr[A|E]$ , etc.

Now we need the probabilities of D, E, and T. What is the probability that all three cards drawn are the same? The probability that the second card is the same as the first is  $\frac{3}{51}$ ; and the probability that the third is the same too is  $\frac{2}{50}$ ; therefore the total probability is  $\frac{\binom{3}{51} \binom{2}{50}}{\binom{3}{51} \binom{2}{50}} = \frac{2}{2550}$ . The probability that all three are unequal is  $\frac{48 \cdot 44}{51 \cdot 50} = \frac{2112}{2550}$ . The probability that two are equal and the third is different is  $3 \frac{3 \cdot 48}{51 \cdot 50} = \frac{432}{2550}$ . Now in half of these cases, the card that is different is higher and in half of the cases it is lower.

Putting this together one gets:

		Uncond. Prob.	Cond. Prob.			Prob. of intersection		
			A	B	C	A	B	C
E	all 3 equal	6/2550	1	0	0	6/2550	0	0
H	2 of 3 equal, 3rd higher	216/2550	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	72/2550	72/2550	72/2550
L	2 of 3 equal, 3rd lower	216/2550	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0	144/2550	72/2550
D	all 3 unequal	2112/2550	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	704/2550	704/2550	704/2550
Sum		2550/2550				926/2550	848/2550	776/2550

I.e., the probability that A wins is  $926/2550 = 463/1275 = .363$ , the probability that B wins is  $848/2550 = 424/1275 = .3325$ , and the probability that C wins is  $776/2550 = 388/1275 = .304$ .

Here we are using  $\Pr[A] = \Pr[A|E] \Pr[E] + \Pr[A|H] \Pr[H] + \Pr[A|L] \Pr[L] + \Pr[A|D] \Pr[D]$ .

PROBLEM 31. *4 points You are the contestant in a game show. There are three closed doors at the back of the stage. Behind one of the doors is a sports car, behind the other two doors are goats. The game master knows which door has the sports car behind it, but you don’t. You have to choose one of the doors; if it is the door with the sports car, the car is yours.*

After you make your choice, say door A, the game master says: “I want to show you something.” He opens one of the two other doors, let us assume it is door B, and it has a goat behind it. Then the game master asks: “Do you still insist on door A, or do you want to reconsider your choice?”

Can you improve your odds of winning by abandoning your previous choice and instead selecting the door which the game master did not open? If so, by how much?

ANSWER. If you switch, you will lose the car if you had initially picked the right door, but you will get the car if you were wrong before! Therefore you improve your chances of winning from  $\frac{1}{3}$  to  $\frac{2}{3}$ . This is simulated on the web, see [www.stat.sc.edu/~west/javahtml/LetsMakeaDeal.htm](http://www.stat.sc.edu/~west/javahtml/LetsMakeaDeal.htm)

It is counterintuitive. You may think that one of the two other doors always has a goat behind it, whatever your choice, therefore there is no reason to switch. But the game master not only shows you that there is another door with a goat, he also shows you one of the other doors with a goat behind it, i.e., he restricts your choice if you switch. This is valuable information. It is as if you could bet on both other doors simultaneously, i.e., you get the car if it is behind one of the doors or C. I.e., if the quiz master had said: I give you the opportunity to switch to the following:

get the car if it is behind  $B$  or  $C$ . Do you want to switch? The only doubt the contestant may have about this is: had I not picked a door with the car behind it then I would not have been offered this opportunity to switch.

□

## 2.8. Ratio of Probabilities as Strength of Evidence

$\Pr_1$  and  $\Pr_2$  are two probability measures defined on the same set  $\mathcal{F}$  of events. Hypothesis  $H_1$  says  $\Pr_1$  is the true probability, and  $H_2$  says  $\Pr_2$  is the true probability. Then the observation of an event  $A$  for which  $\Pr_1[A] > \Pr_2[A]$  is evidence in favor of  $H_1$  as opposed to  $H_2$ . [Roy97] argues that the *ratio* of the probabilities (also called “likelihood ratio”) is the right way to measure the *strength* of this evidence. Among others, the following justification is given [Roy97, p. 7]: If  $H_2$  is true, it is usually not *impossible* to find evidence favoring  $H_1$ , but it is *unlikely*; and its probability is bounded by the (reverse of) the ratio of probabilities.

This can be formulated mathematically as follows: Let  $S$  be the union of all events  $A$  for which  $\Pr_1[A] \geq k \Pr_2[A]$ . Then it can be shown that  $\Pr_2[S] \leq 1/k$ , i.e., if  $H_2$  is true, the probability to find evidence favoring  $H_1$  with strength  $k$  is never greater than  $1/k$ . Here is a proof in the case that there is only a finite number of possible outcomes  $U = \{\omega_1, \dots, \omega_n\}$ : Renumber the outcomes such that for  $i = 1, \dots, m$ ,  $\Pr_1[\{\omega_i\}] < k \Pr_2[\{\omega_i\}]$ , and for  $j = m + 1, \dots, n$ ,  $\Pr_1[\{\omega_j\}] \geq k \Pr_2[\{\omega_j\}]$ . Then  $S = \{\omega_{m+1}, \dots, \omega_n\}$ , therefore  $\Pr_2[S] = \sum_{j=m+1}^n \Pr_2[\{\omega_j\}] \leq \sum_{j=m+1}^n \frac{\Pr_1[\{\omega_j\}]}{k} = \frac{1}{k} \Pr_1[S] \leq \frac{1}{k}$  as claimed. The last inequality holds because  $\Pr_1[S] \leq 1$ , and the equal-sign before this is simply the definition of  $S$ .

With more mathematical effort, see [Rob70], one can strengthen this simple inequality in a very satisfactory manner: Assume an unscrupulous researcher attempts to find evidence supporting his favorite but erroneous hypothesis  $H_1$  over his rival’s  $H_2$  by a factor of at least  $k$ . He proceeds as follows: he observes an outcome of the above experiment once, say the outcome is  $\omega_{i(1)}$ . If  $\Pr_1[\{\omega_{i(1)}\}] \geq k \Pr_2[\{\omega_{i(1)}\}]$  he publishes his result; if not, he makes a second independent observation of the experiment  $\omega_{i(2)}$ . If  $\Pr_1[\{\omega_{i(1)}\}] \Pr_1[\{\omega_{i(2)}\}] > k \Pr_2[\{\omega_{i(1)}\}] \Pr_2[\{\omega_{i(2)}\}]$  he publishes his result; if not he makes a third observation and incorporates that in his publication as well, etc. It can be shown that this strategy will not help: if his rival’s hypothesis is true, then the probability that he will ever be able to publish results which seem to show that his own hypothesis is true is still  $\leq 1/k$ . I.e., the sequence of independent observations  $\omega_{i(2)}, \omega_{i(2)}, \dots$  is such that

$$(2.8.1) \quad \Pr_2 \left[ \prod_{j=1}^n \Pr_1[\{\omega_{i(j)}\}] \geq k \prod_{j=1}^n \Pr_2[\{\omega_{i(j)}\}] \text{ for some } n = 1, 2, \dots \right] \leq \frac{1}{k}$$

It is not possible to take advantage of the indeterminacy of a random outcome by carrying on until chance places one ahead, and then to quit. If one fully discloses

all the evidence one is accumulating, then the probability that this accumulated evidence supports one’s hypothesis cannot rise above  $1/k$ .

**PROBLEM 32.** *It is usually not possible to assign probabilities to the hypotheses  $H_1$  and  $H_2$ , but sometimes it is. Show that in this case, the likelihood ratio of evidence  $A$  is the factor by which the ratio of the probabilities of  $H_1$  and  $H_2$  is changed by observation of  $A$ , i.e.,*

$$(2.8.2) \quad \frac{\Pr[H_1|A]}{\Pr[H_2|A]} = \frac{\Pr[H_1]}{\Pr[H_2]} \frac{\Pr[A|H_1]}{\Pr[A|H_2]}$$

**ANSWER.** Apply Bayes’s theorem (2.9.1) twice, once for the numerator, once for the denominator.

A world in which probability theory applies is therefore a world in which the transitive dimension must be distinguished from the intransitive dimension. Research results are not determined by the goals of the researcher.

## 2.9. Bayes Theorem

In its simplest form Bayes’s theorem reads

$$(2.9.1) \quad \Pr[A|B] = \frac{\Pr[B|A] \Pr[A]}{\Pr[B|A] \Pr[A] + \Pr[B|A'] \Pr[A']}$$

**PROBLEM 33.** *Prove Bayes theorem!*

**ANSWER.** Obvious since numerator is  $\Pr[B \cap A]$  and denominator  $\Pr[B \cap A] + \Pr[B \cap A'] = \Pr[B]$ .

This theorem has its significance in cases in which  $A$  can be interpreted as a cause of  $B$ , and  $B$  an effect of  $A$ . For instance,  $A$  is the event that a student who was picked randomly from a class has learned for a certain exam, and  $B$  is the event that he passed the exam. Then the righthand side expression contains the information which you would know from the cause-effect relations: the unconditional probability of the event which is the cause, and the conditional probabilities of the effect conditioned on whether or not the cause happened. From this, the formula computes the conditional probability of the cause given that the effect happened. Bayes’s theorem tells us therefore: if we know that the effect happened, how sure can we be that the cause happened? Clearly, Bayes’s theorem has relevance to statistical inference.

Let’s stay with the example with learning for the exam; assume  $\Pr[A] = 60\%$ ,  $\Pr[B|A] = .8$ , and  $\Pr[B|A'] = .5$ . Then the probability that a student who passed the exam has learned for it is  $\frac{(.8)(.6)}{(.8)(.6) + (.5)(.4)} = \frac{.48}{.68} = .706$ . Look at these numbers. The numerator is the average percentage of students who learned and passed, and the denominator average percentage of students who passed.

PROBLEM 34. *AIDS diagnostic tests are usually over 99.9% accurate on those who do not have AIDS (i.e., only 0.1% false positives) and 100% accurate on those who have AIDS (i.e., no false negatives at all). (A test is called positive if it indicates that the subject has AIDS.)*

• a. *3 points Assuming that 0.5% of the population actually have AIDS, compute the probability that a particular individual has AIDS, given that he or she has tested positive.*

ANSWER. A is the event that he or she has AIDS, and T the event that the test is positive.

$$\begin{aligned}\Pr[A|T] &= \frac{\Pr[T|A]\Pr[A]}{\Pr[T|A]\Pr[A] + \Pr[T|A']\Pr[A']} = \frac{1 \cdot 0.005}{1 \cdot 0.005 + 0.001 \cdot 0.995} = \\ &= \frac{100 \cdot 0.5}{100 \cdot 0.5 + 0.1 \cdot 99.5} = \frac{1000 \cdot 5}{1000 \cdot 5 + 1 \cdot 995} = \frac{5000}{5995} = \frac{1000}{1199} \approx 0.834028\end{aligned}$$

Even after testing positive there is still a 16.6% chance that this person does not have AIDS.  $\square$

• b. *1 point If one is young, healthy and not in one of the risk groups, then the chances of having AIDS are not 0.5% but 0.1% (this is the proportion of the applicants to the military who have AIDS). Re-compute the probability with this alternative number.*

ANSWER.

$$\frac{1 \cdot 0.001}{1 \cdot 0.001 + 0.001 \cdot 0.999} = \frac{100 \cdot 0.1}{100 \cdot 0.1 + 0.1 \cdot 99.9} = \frac{1000 \cdot 1}{1000 \cdot 1 + 1 \cdot 999} = \frac{1000}{1000 + 999} = \frac{1000}{1999} \approx 0.50025. \quad \square$$

## 2.10. Independence of Events

**2.10.1. Definition of Independence.** Heuristically, we want to say: event B is independent of event A if  $\Pr[B|A] = \Pr[B|A']$ . From this follows by Problem 23 that the conditional probability is equal to the unconditional probability  $\Pr[B]$ , i.e.,  $\Pr[B] = \Pr[B \cap A]/\Pr[A]$ . Therefore we will adopt as definition of independence the so-called *multiplication rule*:

Definition: B and A are independent, notation  $B \perp A$ , if  $\Pr[B \cap A] = \Pr[B]\Pr[A]$ .

This is a symmetric condition, i.e., if B is independent of A, then A is also independent of B. This symmetry is not immediately obvious given the above definition of independence, and it also has the following nontrivial practical implication (this example from [Daw79a, pp. 2/3]): A is the event that one is exposed to some possibly carcinogenic agent, and B the event that one develops a certain kind of cancer. In order to test whether  $B \perp A$ , i.e., whether the exposure to the agent does not increase the incidence of cancer, one often collects two groups of subjects, one group which has cancer and one control group which does not, and checks whether the exposure in these two groups to the carcinogenic agent is the same. I.e., the experiment checks

whether  $A \perp B$ , although the purpose of the experiment was to determine whether  $B \perp A$ .

PROBLEM 35. *3 points Given that  $\Pr[B \cap A] = \Pr[B] \cdot \Pr[A]$  (i.e., B is independent of A), show that  $\Pr[B \cap A'] = \Pr[B] \cdot \Pr[A']$  (i.e., B is also independent of A').*

ANSWER. If one uses our heuristic definition of independence, i.e., B is independent of event A if  $\Pr[B|A] = \Pr[B|A']$ , then it is immediately obvious since definition is symmetric in A and A'. However if we use the multiplication rule as the definition of independence, as the text of this Problem suggests, we have to do a little more work: Since B is the disjoint union of  $(B \cap A)$  and  $(B \cap A')$ , it follows  $\Pr[B] = \Pr[B \cap A] + \Pr[B \cap A']$  or  $\Pr[B \cap A'] = \Pr[B] - \Pr[B \cap A]$ .  $\Pr[B \cap A] = \Pr[B]\Pr[A]$  and  $\Pr[B \cap A'] = \Pr[B] - \Pr[B]\Pr[A] = \Pr[B](1 - \Pr[A]) = \Pr[B]\Pr[A']$ .

PROBLEM 36. *2 points A and B are two independent events with  $\Pr[A] = \frac{1}{3}$  and  $\Pr[B] = \frac{1}{4}$ . Compute  $\Pr[A \cup B]$ .*

ANSWER.  $\Pr[A \cup B] = \Pr[A] + \Pr[B] - \Pr[A \cap B] = \Pr[A] + \Pr[B] - \Pr[A]\Pr[B] = \frac{1}{3} + \frac{1}{4} - \frac{1}{12} = \frac{1}{2}$ .

PROBLEM 37. *3 points You have an urn with five white and five red balls. You take two balls out without replacement. A is the event that the first ball is white and B that the second ball is white. a. What is the probability that the first ball is white? b. What is the probability that the second ball is white? c. What is the probability that both have the same color? d. Are these two events independent, i.e., is  $\Pr[B|A] = \Pr[B]$ ? e. Are these two events disjoint, i.e., is  $A \cap B = \emptyset$ ?*

ANSWER. Clearly,  $\Pr[A] = 1/2$ .  $\Pr[B] = \Pr[B|A]\Pr[A] + \Pr[B|A']\Pr[A'] = (4/9)(1/2) + (5/9)(1/2) = 1/2$ . The events are not independent:  $\Pr[B|A] = 4/9 \neq \Pr[B]$ , or  $\Pr[A \cap B] = \frac{5}{10} \cdot \frac{4}{9} = 2/9 \neq 1/4$ . They would be independent if the first ball had been replaced. The events are also not disjoint: it is possible that both balls are white.

**2.10.2. Independence of More than Two Events.** If there are more than two events, we must require that all possible intersections of these events, not only the pairwise intersections, follow the above multiplication rule. For instance,

$$(2.10.1) \quad \begin{aligned} & \Pr[A \cap B] = \Pr[A]\Pr[B]; \\ & \Pr[A \cap C] = \Pr[A]\Pr[C]; \\ & \Pr[B \cap C] = \Pr[B]\Pr[C]; \\ & \Pr[A \cap B \cap C] = \Pr[A]\Pr[B]\Pr[C] \end{aligned} \iff \begin{aligned} & \text{A, B, C mutually independent} \end{aligned}$$

This last condition is not implied by the other three. Here is an example. Draw a ball at random from an urn containing four balls numbered 1, 2, 3, 4. Define  $A = \{1, 2\}$ ,  $B = \{2, 4\}$ , and  $C = \{3, 4\}$ . These events are pairwise independent but not mutually independent.

PROBLEM 38. 2 points Flip a coin two times independently and define the following three events:

$$(2.10.2) \quad \begin{aligned} A &= \text{Head in first flip} \\ B &= \text{Head in second flip} \\ C &= \text{Same face in both flips.} \end{aligned}$$

Are these three events pairwise independent? Are they mutually independent?

ANSWER.  $\mathcal{U} = \left\{ \begin{smallmatrix} HH & HT \\ TH & TT \end{smallmatrix} \right\}$ .  $A = \{HH, HT\}$ ,  $B = \{HH, TH\}$ ,  $C = \{HH, TT\}$ .  $\Pr[A] = \frac{1}{2}$ ,  $\Pr[B] = \frac{1}{2}$ ,  $\Pr[C] = \frac{1}{2}$ . They are pairwise independent, but  $\Pr[A \cap B \cap C] = \Pr[\{HH\}] = \frac{1}{4} \neq \Pr[A]\Pr[B]\Pr[C]$ , therefore the events cannot be mutually independent.  $\square$

PROBLEM 39. 3 points  $A$ ,  $B$ , and  $C$  are pairwise independent events whose probabilities are greater than zero and smaller than one, and  $A \cap B \subset C$ . Can those events be mutually independent?

ANSWER. No; from  $A \cap B \subset C$  follows  $A \cap B \cap C = A \cap B$  and therefore  $\Pr[A \cap B \cap C] \neq \Pr[A \cap B]\Pr[C]$  since  $\Pr[C] < 1$  and  $\Pr[A \cap B] > 0$ .  $\square$

If one takes unions, intersections, complements of different mutually independent events, one will still end up with mutually independent events. E.g., if  $A$ ,  $B$ ,  $C$  mutually independent, then  $A'$ ,  $B$ ,  $C$  are mutually independent as well, and  $A \cap B$  independent of  $C$ , and  $A \cup B$  independent of  $C$ , etc. This is not the case if the events are only pairwise independent. In Problem 39,  $A \cap B$  is not independent of  $C$ .

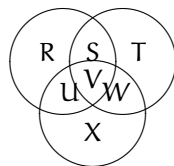


FIGURE 1. Generic Venn Diagram for 3 Events

**2.10.3. Conditional Independence.** If  $A$  and  $B$  are independent in the probability measure conditionally on  $C$ , i.e., if  $\Pr[A \cap B|C] = \Pr[A|C]\Pr[B|C]$ , then they are called conditionally independent given that  $C$  occurred, notation  $A \perp B|C$ . In formulas,

$$(2.10.3) \quad \frac{\Pr[A \cap B|C] \Pr[C]}{\Pr[A|C] \Pr[B|C]} = \frac{\Pr[A \cap B \cap C]}{\Pr[C]}.$$

PROBLEM 40. 5 points Show that  $A \perp B|C$  is equivalent to  $\Pr[A|B \cap C] = \Pr[A|C]$ . In other words: independence of  $A$  and  $B$  conditionally on  $C$  means: once we know that  $C$  occurred, the additional knowledge whether  $B$  occurred or not will not help to sharpen our knowledge about  $A$ .

Literature about conditional independence (of random variables, not of events) includes [Daw79a], [Daw79b], [Daw80].

**2.10.4. Independent Repetition of an Experiment.** If a given experiment has sample space  $\mathcal{U}$ , and we perform the experiment  $n$  times in a row, then the repetition can be considered a single experiment with the sample space consisting of  $n$ -tuples of elements of  $\mathcal{U}$ . This set is called the product set  $\mathcal{U}^n = \mathcal{U} \times \mathcal{U} \times \cdots \times \mathcal{U}$  ( $n$  terms).

If a probability measure  $\Pr$  is given on  $\mathcal{F}$ , then one can define in a unique way a probability measure on the subsets of the product set so that events in different repetitions are always independent of each other.

The *Bernoulli experiment* is the simplest example of such an independent repetition.  $\mathcal{U} = \{s, f\}$  (stands for success and failure). Assume  $\Pr[\{s\}] = p$ , and then the experimenter has several independent trials. For instance,  $\mathcal{U}^5$  has, among other things, the following possible outcomes:

$$(2.10.4) \quad \begin{array}{ll} \text{If } \omega = (f, f, f, f, f) & \text{then } \Pr[\{\omega\}] = (1-p)^n \\ (f, f, f, f, s) & (1-p)^{n-1}p \\ (f, f, f, s, f) & (1-p)^{n-1}p \\ (f, f, f, s, s) & (1-p)^{n-2}p^2 \\ (f, f, s, f, f) & (1-p)^{n-1}p, \text{ etc.} \end{array}$$

One sees, this is very cumbersome, and usually unnecessarily so. If we toss a coin 5 times, the only thing we usually want to know is how many successes there were. As long as the experiments are independent, the question how the successes were distributed over the  $n$  different trials is far less important. This brings us to the definition of a random variable, and to the concept of a sufficient statistic.

## 2.11. How to Plot Frequency Vectors and Probability Vectors

If there are only 3 possible outcomes, i.e.,  $\mathcal{U} = \{\omega_1, \omega_2, \omega_3\}$ , then the set of probability measures is the set of nonnegative 3-vectors whose components sum up to 1. Graphically, such vectors can be represented as points inside a trilateral triangle with height 1: the three components of the vector are the distances of the point to each of the sides of the triangle. The R/Spplus-function `triplot` in the `ecm` package, written by Jim Ramsay `ramsay@ramsay2.psych.mcgill.ca`, does this, with optional rescaling if the rows of the data matrix do not have unit sums.

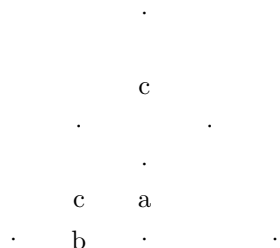


FIGURE 2. Geometry of an equilateral triangle

PROBLEM 41. In an equilateral triangle, call  $a$  = the distance of the sides from the center point,  $b$  = half the side length, and  $c$  = the distance of the corners from the center point (as in Figure 2). Show that  $b = a\sqrt{3}$  and  $c = 2a$ .

ANSWER. From  $(a + c)^2 + b^2 = 4b^2$ , i.e.,  $(a + c)^2 = 3b^2$ , follows  $a + c = b\sqrt{3}$ . But we also have  $a^2 + b^2 = c^2$ . Therefore  $a^2 + 2ac + c^2 = 3b^2 = 3c^2 - 3a^2$ , or  $4a^2 + 2ac - 2c^2 = 0$  or  $2a^2 + ac - c^2 = (2a - c)(a + c) = 0$ . The positive solution is therefore  $c = 2a$ . This gives  $a + c = 3a = b\sqrt{3}$ , or  $b = a\sqrt{3}$ .  $\square$

And the function `quadplot`, also written by Jim Ramsey, does *quadrilinear* plots, meaning that proportions for four categories are plotted within a regular tetrahedron. `Quadplot` displays the probability tetrahedron and its points using `XGobi`. Each vertex of the triangle or tetrahedron corresponds to the degenerate probability distribution in which one of the events has probability 1 and the others have probability 0. The labels of these vertices indicate which event has probability 1.

The script `kai` is an example visualizing data from [Mor65]; it can be run using the command `ecmet.script(kai)`.

Example: Statistical linguistics.

In the study of ancient literature, the authorship of texts is a perplexing problem. When books were written and reproduced by hand, the rights of authorship were limited and what would now be considered forgery was common. The names of reputable authors were borrowed in order to sell books, get attention for books, or the writings of disciples and collaborators were published under the name of the master, or anonymous old manuscripts were optimistically attributed to famous authors. In the absence of conclusive evidence of authorship, the attribution of ancient texts must be based on the texts themselves, for instance, by statistical analysis of literary style. Here it is necessary to find stylistic criteria which vary from author to author, but are independent of the subject matter of the text. An early suggestion was to use the probability distribution of word length, but this was never acted upon, because it is too dependent on the subject matter. Sentence-length distributions, on the other hand, have proved highly reliable. [Mor65, p. 184] says that sentence-length

is “periodic rather than random,” therefore the sample should have at least about 100 sentences. “Sentence-length distributions are not suited to dialogue, they cannot be used on commentaries written on one author by another, nor are they reliable for such texts as the fragmentary books of the historian Diodorus Siculus.”

PROBLEM 42. According to [Mor65, p. 184], sentence-length is “periodic rather than random.” What does this mean?

ANSWER. In a text, passages with long sentences alternate with passages with shorter sentences. This is why one needs at least 100 sentences to get a representative distribution of sentence lengths, and this is why fragments and drafts and commentaries on others’ writings do not exhibit an average sentence length distribution: they do not have the melody of the finished text.

Besides the *length* of sentences, also the number of common words which express a general relation (“and”, “in”, “but”, “I”, “to be”) is random with the same distribution at least among the same genre. By contrast, the occurrence of the definite article “the” cannot be modeled by simple probabilistic laws because the number of nouns with definite article depends on the subject matter.

Table 1 has data about the epistles of St. Paul. Abbreviations: **Rom** Romans; **Co1** 1st Corinthians; **Co2** 2nd Corinthians; **Gal** Galatians; **Phi** Philippians; **Col** Colossians; **Th1** 1st Thessalonians; **Ti1** 1st Timothy; **Ti2** 2nd Timothy; **Heb** Hebrews. 2nd Thessalonians, Titus, and Philemon were excluded because they were too short to give reliable samples. From an analysis of these and other data [Mor65, p. 224] the first 4 epistles (Romans, 1st Corinthians, 2nd Corinthians, and Galatians) form a consistent group, and all the other epistles lie more than 2 standard deviations from the mean of this group (using  $\chi^2$  statistics). If Paul is defined as being the author of Galatians, then he also wrote Romans and 1st and 2nd Corinthians. The remaining epistles come from at least six hands.

TABLE 1. Number of Sentences in Paul’s Epistles with 0, 1, 2, and  $\geq 3$  occurrences of *kai*

	Rom	Co1	Co2	Gal	Phi	Col	Th1	Ti1	Ti2	Heb
no <i>kai</i>	386	424	192	128	42	23	34	49	45	155
one	141	152	86	48	29	32	23	38	28	94
two	34	35	28	5	19	17	8	9	11	37
3 or more	17	16	13	6	12	9	16	10	4	24

PROBLEM 43. Enter the data from Table 1 into `xgobi` and brush the four epistles which are, according to Morton, written by Paul himself. 3 of those points are almost on top of each other, and one is a little apart. Which one is this?

ANSWER. In R, issue the commands `library(xgobi)` then `data(PaulKAI)` then `quadplot(PaulKAI, normalize = TRUE)`. If you have `xgobi` but not R, this dataset is one of the default datasets coming with `xgobi`.

□

## CHAPTER 3

## Random Variables

## 3.1. Notation

Throughout these class notes, lower case bold letters will be used for vectors and upper case bold letters for matrices, and letters that are not bold for scalars. The  $(i, j)$  element of the matrix  $\mathbf{A}$  is  $a_{ij}$ , and the  $i$ th element of a vector  $\mathbf{b}$  is  $b_i$ ; the arithmetic mean of all elements is  $\bar{b}$ . All vectors are column vectors; if a row vector is needed, it will be written in the form  $\mathbf{b}^\top$ . Furthermore, the on-line version of these notes uses green symbols for random variables, and the corresponding black symbols for the values taken by these variables. If a black-and-white printout of the on-line version is made, then the symbols used for random variables and those used for specific values taken by these random variables can only be distinguished by their grey scale or cannot be distinguished at all; therefore a special monochrome version is available which should be used for the black-and-white printouts. It uses an *upright* math font, called “Euler,” for the random variables, and the same letter in the usual slanted italic font for the values of these random variables.

Example: If  $\mathbf{y}$  is a random vector, then  $\mathbf{y}$  denotes a particular value, for instance an observation, of the whole vector;  $y_i$  denotes the  $i$ th element of  $\mathbf{y}$  (a random scalar), and  $y_i$  is a particular value taken by that element (a nonrandom scalar).

With real-valued random variables, the powerful tools of calculus become available to us. Therefore we will begin the chapter about random variables with a digression about infinitesimals

## 3.2. Digression about Infinitesimals

In the following pages we will recapitulate some basic facts from calculus. But it will differ in two respects from the usual calculus classes. (1) everything will be given its probability-theoretic interpretation, and (2) we will make explicit use of infinitesimals. This last point bears some explanation.

You may say infinitesimals do not exist. Do you know the story with Achilles and the turtle? They are racing, the turtle starts 1 km ahead of Achilles, and Achilles runs ten times as fast as the turtle. So when Achilles arrives at the place the turtle started, the turtle has run 100 meters; and when Achilles has run those 100 meters,

the turtle has run 10 meters, and when Achilles has run the 10 meters, then the turtle has run 1 meter, etc. The Greeks were actually arguing whether Achilles would ever reach the turtle.

This may sound like a joke, but in some respects, modern mathematics never went beyond the level of the Greek philosophers. If a modern mathematician says something like

$$(3.2.1) \quad \lim_{i \rightarrow \infty} \frac{1}{i} = 0, \quad \text{or} \quad \lim_{n \rightarrow \infty} \sum_{i=0}^n \frac{1}{10^i} = \frac{10}{9},$$

then he will probably say that the lefthand term in each equation never really reaches the number written on the right, all he will say is that the term on the left comes *arbitrarily close* to it.

This is like saying: I know that Achilles will get as close as 1 cm or 1 mm to the turtle, he will get closer than any distance, however small, to the turtle, instead of simply saying that Achilles reaches the turtle. Modern mathematical proofs are full of races between Achilles and the turtle of the kind: give me an  $\varepsilon$ , and I will prove you that the thing will come at least as close as  $\varepsilon$  to its goal (so-called epsilon-tolerance) but never speaking about the moment when the thing will reach its goal.

Of course, it “works,” but it makes things terribly cumbersome, and it may have prevented people from seeing connections.

Abraham Robinson in [Rob74] is one of the mathematicians who tried to remedy this. He did it by adding more numbers, infinite numbers and infinitesimal numbers. Robinson showed that one can use infinitesimals without getting into contradictions, and he demonstrated that mathematics becomes much more intuitive this way, not only its elementary proofs, but especially the deeper results. One of the elementary books based on his calculus is [HK79].

The well-known logician Kurt Gödel said about Robinson’s work: “I think, in the coming years it will be considered a great oddity in the history of mathematics that the first exact theory of infinitesimals was developed 300 years after the invention of the differential calculus.”

Gödel called Robinson’s theory the *first* theory. I would like to add here the following speculation: perhaps Robinson shares the following error with the “standard” mathematicians whom he criticizes: they consider numbers only in a static way, without allowing them to move. It would be beneficial to expand on the intuition of the inventors of differential calculus, who talked about “fluxions,” i.e., quantities in flux, in motion. Modern mathematicians even use arrows in their symbol for limits, but they are not calculating with moving quantities, only with static quantities.

This perspective makes the category-theoretical approach to infinitesimals taken in [MR91] especially promising. Category theory considers objects on the same footing with their transformations (and uses lots of arrows).

Maybe a few years from now mathematics will be done right. We should not let this temporary backwardness of mathematics allow to hold us back in our *intuition*. The equation  $\frac{\Delta y}{\Delta x} = 2x$  does not hold exactly on a parabola for any pair of given (static)  $\Delta x$  and  $\Delta y$ ; but if you take a pair  $(\Delta x, \Delta y)$  which is *moving* towards zero then this equation holds *in the moment when they reach zero, i.e., when they vanish*. Writing  $dy$  and  $dx$  means therefore: we are looking at magnitudes which are in the process of vanishing. If one applies a function to a moving quantity one again gets a moving quantity, and the derivative of this function compares the speed with which the transformed quantity moves with the speed of the original quantity. Likewise, the equation  $\sum_{i=1}^n \frac{1}{2^n} = 1$  holds *in the moment when  $n$  reaches infinity*. From this point of view, the axiom of  $\sigma$ -additivity in probability theory (in its equivalent form of rising or declining sequences of events) indicates that the probability of a vanishing event vanishes.

Whenever we talk about infinitesimals, therefore, we really mean magnitudes which are moving, and which are in the process of vanishing.  $dV_{x,y}$  is therefore not, as one might think from what will be said below, a static but small volume element located close to the point  $(x, y)$ , but it is a volume element which is vanishing into the point  $(x, y)$ . The probability density function therefore signifies the speed with which the probability of a vanishing element vanishes.

### 3.3. Definition of a Random Variable

The best intuition of a random variable would be to view it as a numerical variable whose values are not determinate but follow a statistical pattern, and call it  $x$ , while possible values of  $x$  are called  $x$ .

In order to make this a mathematically sound definition, one says: A mapping  $x : U \rightarrow \mathbb{R}$  of the set  $U$  of all possible outcomes into the real numbers  $\mathbb{R}$  is called a random variable. (Again, mathematicians are able to construct pathological mappings that cannot be used as random variables, but we let that be their problem, not ours.) The green  $x$  is then defined as  $x = x(\omega)$ . I.e., all the randomness is shunted off into the process of selecting an element of  $U$ . Instead of being an indeterminate function, it is defined as a determinate function of the random  $\omega$ . It is written here as  $x(\omega)$  and not as  $x(\omega)$  because the function itself is determinate, only its argument is random.

Whenever one has a mapping  $x : U \rightarrow \mathbb{R}$  between sets, one can construct from it in a natural way an “inverse image” mapping between subsets of these sets. Let  $\mathcal{F}$ , as usual, denote the set of subsets of  $U$ , and let  $\mathcal{B}$  denote the set of subsets of  $\mathbb{R}$ . We will define a mapping  $x^{-1} : \mathcal{B} \rightarrow \mathcal{F}$  in the following way: For any  $B \subset \mathbb{R}$ , we define  $x^{-1}(B) = \{\omega \in U : x(\omega) \in B\}$ . (This is not the usual inverse of a mapping, which does not always exist. The inverse-image mapping always exists, but the inverse image of a one-element set is no longer necessarily a one-element set; it may have more than one element or may be the empty set.)

This “inverse image” mapping is well behaved with respect to unions and intersections, etc. In other words, we have identities  $x^{-1}(A \cap B) = x^{-1}(A) \cap x^{-1}(B)$  and  $x^{-1}(A \cup B) = x^{-1}(A) \cup x^{-1}(B)$ , etc.

PROBLEM 44. *Prove the above two identities.*

ANSWER. These are a very subtle proofs.  $x^{-1}(A \cap B) = \{\omega \in U : x(\omega) \in A \cap B\} = \{\omega \in U : x(\omega) \in A \text{ and } x(\omega) \in B\} = \{\omega \in U : x(\omega) \in A\} \cap \{\omega \in U : x(\omega) \in B\} = x^{-1}(A) \cap x^{-1}(B)$ . The other identity has a similar proof.

PROBLEM 45. *Show, on the other hand, by a counterexample, that the “direct image” mapping defined by  $x(E) = \{r \in \mathbb{R} : \text{there exists } \omega \in E \text{ with } x(\omega) = r\}$  no longer satisfies  $x(E \cap F) = x(E) \cap x(F)$ .*

By taking inverse images under a random variable  $x$ , the probability measure on  $\mathcal{F}$  is transplanted into a probability measure on the subsets of  $\mathbb{R}$  by the simple prescription  $\Pr[B] = \Pr[x^{-1}(B)]$ . Here,  $B$  is a subset of  $\mathbb{R}$  and  $x^{-1}(B)$  one of  $\mathcal{F}$ ,  $\Pr$  on the right side is the given probability measure on  $\mathcal{F}$ , while the  $\Pr$  on the left side is the new probability measure on  $\mathbb{R}$  induced by  $x$ . This induced probability measure is called the probability law or probability distribution of the random variable.

Every random variable induces therefore a probability measure on  $\mathbb{R}$ , and this probability measure, not the mapping itself, is the most important ingredient of a random variable. That is why Amemiya’s first definition of a random variable (definition 3.1.1 on p. 18) is: “A random variable is a variable that takes values according to a certain distribution.” In other words, it is the outcome of an experiment whose set of possible outcomes is  $\mathbb{R}$ .

### 3.4. Characterization of Random Variables

We will begin our systematic investigation of random variables with an overview over all possible probability measures on  $\mathbb{R}$ .

The simplest way to get such an overview is to look at the *cumulative distribution functions*. Every probability measure on  $\mathbb{R}$  has a cumulative distribution function, but we will follow the common usage of assigning the cumulative distribution function to a probability measure but to the random variable which induces this probability measure on  $\mathbb{R}$ .

Given a random variable  $x : U \ni \omega \mapsto x(\omega) \in \mathbb{R}$ . Then the cumulative distribution function of  $x$  is the function  $F_x : \mathbb{R} \rightarrow \mathbb{R}$  defined by:

$$(3.4.1) \quad F_x(a) = \Pr[\{\omega \in U : x(\omega) \leq a\}] = \Pr[x \leq a].$$

This function uniquely defines the probability measure which  $x$  induces on  $\mathbb{R}$ .



Properties of cumulative distribution functions: a function  $F: \mathbb{R} \rightarrow \mathbb{R}$  is a cumulative distribution function if and only if

$$(3.4.2) \quad a \leq b \Rightarrow F(a) \leq F(b)$$

$$(3.4.3) \quad \lim_{a \rightarrow -\infty} F(a) = 0$$

$$(3.4.4) \quad \lim_{a \rightarrow \infty} F(a) = 1$$

$$(3.4.5) \quad \lim_{\varepsilon \rightarrow 0, \varepsilon > 0} F(a + \varepsilon) = F(a)$$

Equation (3.4.5) is the definition of *continuity from the right* (because the limit holds only for  $\varepsilon \geq 0$ ). Why is a cumulative distribution function continuous from the right? For every nonnegative sequence  $\varepsilon_1, \varepsilon_2, \dots \geq 0$  converging to zero which also satisfies  $\varepsilon_1 \geq \varepsilon_2 \geq \dots$  follows  $\{x \leq a\} = \bigcap_i \{x \leq a + \varepsilon_i\}$ ; for these sequences, therefore, the statement follows from what Problem 14 above said about the probability of the intersection of a declining set sequence. And a converging sequence of nonnegative  $\varepsilon_i$  which is not declining has a declining subsequence.

A cumulative distribution function need not be continuous from the left. If  $\lim_{\varepsilon \rightarrow 0, \varepsilon > 0} F(x - \varepsilon) \neq F(x)$ , then  $x$  is a jump point, and the height of the jump is the probability that  $x = x$ .

It is a matter of convention whether we are working with right continuous or left continuous functions here. If the distribution function were defined as  $\Pr[x < a]$  (some authors do this, compare [Ame94, p. 43]), then it would be continuous from the left but not from the right.

**PROBLEM 46.** 6 points Assume  $F_x(x)$  is the cumulative distribution function of the random variable  $x$  (whose distribution is not necessarily continuous). Which of the following formulas are correct? Give proofs or verbal justifications.

$$(3.4.6) \quad \Pr[x = x] = \lim_{\varepsilon > 0; \varepsilon \rightarrow 0} F_x(x + \varepsilon) - F_x(x)$$

$$(3.4.7) \quad \Pr[x = x] = F_x(x) - \lim_{\delta > 0; \delta \rightarrow 0} F_x(x - \delta)$$

$$(3.4.8) \quad \Pr[x = x] = \lim_{\varepsilon > 0; \varepsilon \rightarrow 0} F_x(x + \varepsilon) - \lim_{\delta > 0; \delta \rightarrow 0} F_x(x - \delta)$$

**ANSWER.** (3.4.6) does not hold generally, since its rhs is always = 0; the other two equations always hold.  $\square$

**PROBLEM 47.** 4 points Assume the distribution of  $z$  is symmetric about zero, i.e.,  $\Pr[z < -z] = \Pr[z > z]$  for all  $z$ . Call its cumulative distribution function  $F_z(z)$ . Show that the cumulative distribution function of the random variable  $q = z^2$  is  $F_q(q) = 2F_z(\sqrt{q}) - 1$  for  $q \geq 0$ , and 0 for  $q < 0$ .

**ANSWER.** If  $q \geq 0$  then

$$(3.4.9) \quad F_q(q) = \Pr[z^2 \leq q] = \Pr[-\sqrt{q} \leq z \leq \sqrt{q}]$$

$$(3.4.10) \quad = \Pr[z \leq \sqrt{q}] - \Pr[z < -\sqrt{q}]$$

$$(3.4.11) \quad = \Pr[z \leq \sqrt{q}] - \Pr[z > \sqrt{q}]$$

$$(3.4.12) \quad = F_z(\sqrt{q}) - (1 - F_z(\sqrt{q}))$$

$$(3.4.13) \quad = 2F_z(\sqrt{q}) - 1.$$

Instead of the cumulative distribution function  $F_y$  one can also use the *quantile function*  $F_y^{-1}$  to characterize a probability measure. As the notation suggests the quantile function can be considered some kind of “inverse” of the cumulative distribution function. The quantile function is the function  $(0, 1) \rightarrow \mathbb{R}$  defined by

$$(3.4.14) \quad F_y^{-1}(p) = \inf\{u : F_y(u) \geq p\}$$

or, plugging the definition of  $F_y$  into (3.4.14),

$$(3.4.15) \quad F_y^{-1}(p) = \inf\{u : \Pr[y \leq u] \geq p\}.$$

The quantile function is only defined on the open unit interval, not on the endpoints 0 and 1, because it would often assume the values  $-\infty$  and  $+\infty$  on these endpoints and the information given by these values is redundant. The quantile function is continuous from the left, i.e., from the other side than the cumulative distribution function. If  $F$  is continuous and strictly increasing, then the quantile function is the inverse of the distribution function in the usual sense, i.e.,  $F^{-1}(F(t)) = t$  for all  $t \in \mathbb{R}$ , and  $F(F^{-1}(p)) = p$  for all  $p \in (0, 1)$ . But even if  $F$  is flat on certain intervals, and/or  $F$  has jump points, i.e.,  $F$  does not have an inverse function, the following important identity holds for every  $y \in \mathbb{R}$  and  $p \in (0, 1)$ :

$$(3.4.16) \quad p \leq F_y(y) \quad \text{iff} \quad F_y^{-1}(p) \leq y$$

**PROBLEM 48.** 3 points Prove equation (3.4.16).

**ANSWER.**  $\Rightarrow$  is trivial: if  $F(y) \geq p$  then of course  $y \geq \inf\{u : F(u) \geq p\}$ .  $\Leftarrow$ :  $y \geq \inf\{u : F(u) \geq p\}$  means that every  $z > y$  satisfies  $F(z) \geq p$ ; therefore, since  $F$  is continuous from the right, also  $F(y) \geq p$ . This proof is from [Rei89, p. 318].

**PROBLEM 49.** You throw a pair of dice and your random variable  $x$  is the sum of the points shown.

- a. Draw the cumulative distribution function of  $x$ .

**ANSWER.** This is Figure 1: the cdf is 0 in  $(-\infty, 2)$ ,  $1/36$  in  $[2,3)$ ,  $3/36$  in  $[3,4)$ ,  $6/36$  in  $[4,5)$ ,  $10/36$  in  $[5,6)$ ,  $15/36$  in  $[6,7)$ ,  $21/36$  in  $[7,8)$ ,  $26/36$  in  $[8,9)$ ,  $30/36$  in  $[9,10)$ ,  $33/36$  in  $[10,11)$ ,  $35/36$  in  $[11,12)$ , and 1 in  $[12, +\infty)$ .

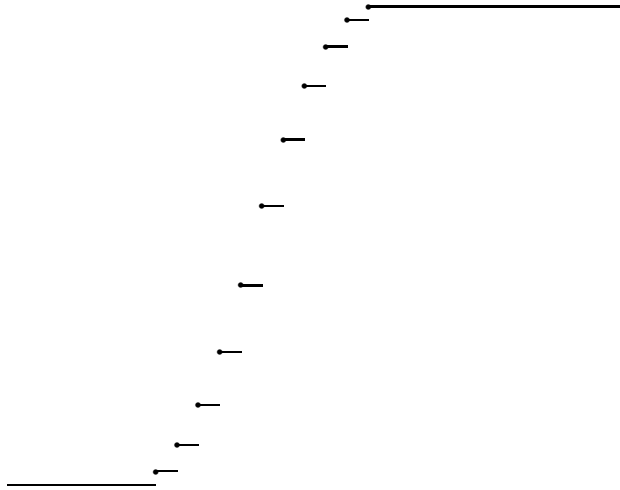


FIGURE 1. Cumulative Distribution Function of Discrete Variable

- b. Draw the quantile function of  $x$ .

ANSWER. This is Figure 2: the quantile function is 2 in  $(0, 1/36]$ , 3 in  $(1/36, 3/36]$ , 4 in  $(3/36, 6/36]$ , 5 in  $(6/36, 10/36]$ , 6 in  $(10/36, 15/36]$ , 7 in  $(15/36, 21/36]$ , 8 in  $(21/36, 26/36]$ , 9 in  $(26/36, 30/36]$ , 10 in  $(30/36, 33/36]$ , 11 in  $(33/36, 35/36]$ , and 12 in  $(35/36, 1]$ .  $\square$

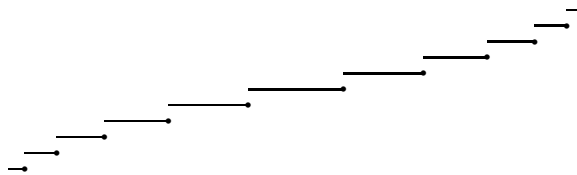


FIGURE 2. Quantile Function of Discrete Variable

PROBLEM 50. 1 point Give the formula of the cumulative distribution function of a random variable which is uniformly distributed between 0 and  $b$ .

ANSWER. 0 for  $x \leq 0$ ,  $x/b$  for  $0 \leq x \leq b$ , and 1 for  $x \geq b$ .  $\square$

### Empirical Cumulative Distribution Function:

Besides the cumulative distribution function of a random variable or of a probability measure, one can also define the empirical cumulative distribution function of a sample. Empirical cumulative distribution functions are zero for all values below the lowest observation, then  $1/n$  for everything below the second lowest, etc. They are step functions. If two observations assume the same value, then the step at that value is twice as high, etc. The empirical cumulative distribution function can be considered an estimate of the cumulative distribution function of the probability distribution underlying the sample. [Rei89, p. 12] writes it as a sum of indicator functions:

$$(3.4.17) \quad F = \frac{1}{n} \sum_i 1_{[x_i, +\infty)}$$

### 3.5. Discrete and Absolutely Continuous Probability Measures

One can define two main classes of probability measures on  $\mathbb{R}$ :

One kind is concentrated in countably many points. Its probability distribution can be defined in terms of the probability mass function.

PROBLEM 51. Show that a distribution function can only have countably many jump points.

ANSWER. Proof: There are at most two with jump height  $\geq \frac{1}{2}$ , at most four with jump height  $\geq \frac{1}{4}$ , etc.

Among the other probability measures we are only interested in those which can be represented by a density function (absolutely continuous). A density function is a nonnegative integrable function which, integrated over the whole line, gives 1. Given such a density function, called  $f_x(x)$ , the probability  $\Pr[x \in (a, b)] = \int_a^b f_x(x) dx$ . The density function is therefore an alternate way to characterize a probability measure. But not all probability measures have density functions.

Those who are not familiar with integrals should read up on them at this point. Start with derivatives, then: the indefinite integral of a function is a function whose derivative is the given function. Then it is an important theorem that the area under the curve is the difference of the values of the indefinite integral at the end points. This is called the definite integral. (The area is considered negative when the curve is below the  $x$ -axis.)

The intuition of a density function comes out more clearly in terms of infinitesimals. If  $f_x(x)$  is the value of the density function at the point  $x$ , then the probability that the outcome of  $x$  lies in an interval of infinitesimal length located near the point  $x$  is the length of this interval, multiplied by  $f_x(x)$ . In formulas, for an infinitesimal

$dx$  follows

$$(3.5.1) \quad \Pr[x \in [x, x + dx]] = f_x(x) |dx|.$$

The name “density function” is therefore appropriate: it indicates how densely the probability is spread out over the line. It is, so to say, the quotient between the probability measure induced by the variable, and the length measure on the real numbers.

If the cumulative distribution function has everywhere a derivative, this derivative is the density function.

### 3.6. Transformation of a Scalar Density Function

Assume  $x$  is a random variable with values in the region  $A \subset \mathbb{R}$ , i.e.,  $\Pr[x \notin A] = 0$ , and  $t$  is a one-to-one mapping  $A \rightarrow \mathbb{R}$ . One-to-one (as opposed to many-to-one) means: if  $a, b \in A$  and  $t(a) = t(b)$ , then already  $a = b$ . We also assume that  $t$  has a continuous nonnegative first derivative  $t' \geq 0$  everywhere in  $A$ . Define the random variable  $y$  by  $y = t(x)$ . We know the density function of  $y$ , and we want to get that of  $x$ . (I.e.,  $t$  expresses the old variable, that whose density function we know, in terms of the new variable, whose density function we want to know.)

Since  $t$  is one-to-one, it follows for all  $a, b \in A$  that  $a = b \iff t(a) = t(b)$ . And recall the definition of a derivative in terms of infinitesimals  $dx$ :  $t'(x) = \frac{t(x+dx) - t(x)}{dx}$ .

In order to compute  $f_x(x)$  we will use the following identities valid for all  $x \in A$ :

$$(3.6.1) \quad f_x(x) |dx| = \Pr[x \in [x, x + dx]] = \Pr[t(x) \in [t(x), t(x + dx)]]$$

$$(3.6.2) \quad = \Pr[t(x) \in [t(x), t(x) + t'(x) dx]] = f_y(t(x)) |t'(x) dx|$$

Absolute values are multiplicative, i.e.,  $|t'(x) dx| = |t'(x)| |dx|$ ; divide by  $|dx|$  to get

$$(3.6.3) \quad f_x(x) = f_y(t(x)) |t'(x)|.$$

This is the transformation formula how to get the density of  $x$  from that of  $y$ . This formula is valid for all  $x \in A$ ; the density of  $x$  is 0 for all  $x \notin A$ .

Heuristically one can get this transformation as follows: write  $|t'(x)| = \frac{|dy|}{|dx|}$ , then one gets it from  $f_x(x) |dx| = f_y(t(x)) |dy|$  by just dividing both sides by  $|dx|$ .

In other words, this transformation rule consists of 4 steps: (1) Determine  $A$ , the range of the new variable; (2) obtain the transformation  $t$  which expresses the old variable in terms of the new variable, and check that it is one-to-one on  $A$ ; (3) plug expression (2) into the old density; (4) multiply this plugged-in density by the absolute value of the derivative of expression (2). This gives the density inside  $A$ ; it is 0 outside  $A$ .

An alternative proof is conceptually simpler but cannot be generalized to the multivariate case: First assume  $t$  is monotonically *increasing*. Then  $F_x(x) = \Pr[x \leq$

$x] = \Pr[t(x) \leq t(i)] = F_y(t(x))$ . Now differentiate and use the chain rule. This also do the monotonically *decreasing* case. This is how [Ame94, theorem 3.6.1 pp. 48] does it. [Ame94, pp. 52/3] has an extension of this formula to many-to-one functions.

PROBLEM 52. 4 points [Lar82, example 3.5.4 on p. 148] Suppose  $y$  has density function

$$(3.6.4) \quad f_y(y) = \begin{cases} 1 & \text{for } 0 < y < 1 \\ 0 & \text{otherwise.} \end{cases}$$

Obtain the density  $f_x(x)$  of the random variable  $x = -\log y$ .

ANSWER. (1) Since  $y$  takes values only between 0 and 1, its logarithm takes values between  $-\infty$  and 0, the negative logarithm therefore takes values between 0 and  $+\infty$ , i.e.,  $A = \{x : 0 < x < +\infty\}$ . (2) Express  $y$  in terms of  $x$ :  $y = e^{-x}$ . This is one-to-one on the whole line, therefore also on  $A$ . (3) Plugging  $y = e^{-x}$  into the density function gives the number 1, since the density function does not depend on the precise value of  $y$ , as long as we know that  $0 < y < 1$  (which we do). (4) The derivative of  $y = e^{-x}$  is  $-e^{-x}$ . As a last step one has to multiply the number 1 by the absolute value of the derivative to get the density inside  $A$ . Therefore  $f_x(x) = e^{-x}$  for  $x > 0$  and 0 otherwise.

PROBLEM 53. 6 points [Dhr86, p. 1574] Assume the random variable  $z$  has the exponential distribution with parameter  $\lambda$ , i.e., its density function is  $f_z(z) = \lambda \exp(-\lambda z)$  for  $z > 0$  and 0 for  $z \leq 0$ . Define  $u = -\log z$ . Show that the density function of  $u$  is  $f_u(u) = \exp(\mu - u - \exp(\mu - u))$  where  $\mu = \log \lambda$ . This density function can be used in Problem 140.

ANSWER. (1) Since  $z$  only has values in  $(0, \infty)$ , its log is well defined, and  $A = \mathbb{R}$ . (2) Express old variable in terms of new:  $-u = \log z$  therefore  $z = e^{-u}$ ; this is one-to-one everywhere. (3) plugging in (since  $e^{-u} > 0$  for all  $u$ , we must plug it into  $\lambda \exp(-\lambda z)$ ) gives  $\lambda \exp(-\lambda e^{-u})$ . (4) the derivative of  $z = e^{-u}$  is  $-e^{-u}$ , taking absolute values gives the Jacobian factor  $e^{-u}$ . Plugging in and multiplying gives the density of  $u$ :  $f_u(u) = \lambda \exp(-\lambda e^{-u}) e^{-u} = \lambda e^{-u - \lambda e^{-u}}$ , and using  $\lambda \exp(-u) = \exp(\mu - u)$  this simplifies to the formula above.

Alternative without transformation rule for densities:  $F_u(u) = \Pr[u \leq u] = \Pr[-\log z \leq u] = \Pr[\log z \geq -u] = \Pr[z \geq e^{-u}] = \int_{e^{-u}}^{+\infty} \lambda e^{-\lambda z} dz = -e^{-\lambda z} \Big|_{e^{-u}}^{+\infty} = e^{-\lambda e^{-u}}$ , now differentiate.

PROBLEM 54. 4 points Assume the random variable  $z$  has the exponential distribution with  $\lambda = 1$ , i.e., its density function is  $f_z(z) = \exp(-z)$  for  $z \geq 0$  and 0 for  $z < 0$ . Define  $u = \sqrt{z}$ . Compute the density function of  $u$ .

ANSWER. (1)  $A = \{u : u \geq 0\}$  since  $\sqrt{\cdot}$  always denotes the nonnegative square root; (2) Express old variable in terms of new:  $z = u^2$ , this is one-to-one on  $A$  (but not one-to-one on all of  $\mathbb{R}$ ); (3) then the derivative is  $2u$ , which is nonnegative as well, no absolute values are necessary; (4) multiplying gives the density of  $u$ :  $f_u(u) = 2u \exp(-u^2)$  if  $u \geq 0$  and 0 elsewhere.

### 3.7. Example: Binomial Variable

Go back to our Bernoulli trial with parameters  $p$  and  $n$ , and define a random variable  $x$  which represents the number of successes. Then the probability mass function of  $x$  is

$$(3.7.1) \quad p_x(k) = \Pr[x=k] = \binom{n}{k} p^k (1-p)^{(n-k)} \quad k = 0, 1, 2, \dots, n$$

Proof is simple, every subset of  $k$  elements represents one possibility of spreading out the  $k$  successes.

We will call any observed random variable a *statistic*. And we call a statistic  $t$  *sufficient for a parameter*  $\theta$  if and only if for any event  $A$  and for any possible value  $t$  of  $t$ , the conditional probability  $\Pr[A|t \leq t]$  does not involve  $\theta$ . This means: after observing  $t$  no additional information can be obtained about  $\theta$  from the outcome of the experiment.

**PROBLEM 55.** *Show that  $x$ , the number of successes in the Bernoulli trial with parameters  $p$  and  $n$ , is a sufficient statistic for the parameter  $p$  (the probability of success), with  $n$ , the number of trials, a known fixed number.*

**ANSWER.** Since the distribution of  $x$  is discrete, it is sufficient to show that for any given  $k$ ,  $\Pr[A|x=k]$  does not involve  $p$  whatever the event  $A$  in the Bernoulli trial. Furthermore, since the Bernoulli trial with  $n$  tries is finite, we only have to show it if  $A$  is an elementary event in  $\mathcal{F}$ , i.e., an event consisting of one element. Such an elementary event would be that the outcome of the trial has a certain given sequence of successes and failures. A general  $A$  is the finite disjoint union of all elementary events contained in it, and if the probability of each of these elementary events does not depend on  $p$ , then their sum does not either.

Now start with the definition of conditional probability

$$(3.7.2) \quad \Pr[A|x=k] = \frac{\Pr[A \cap \{x=k\}]}{\Pr[x=k]}.$$

If  $A$  is an elementary event whose number of successes is not  $k$ , then  $A \cap \{x=k\} = \emptyset$ , therefore its probability is 0, which does not involve  $p$ . If  $A$  is an elementary event which has  $k$  successes, then  $A \cap \{x=k\} = A$ , which has probability  $p^k(1-p)^{n-k}$ . Since  $\Pr[\{x=k\}] = \binom{n}{k} p^k(1-p)^{n-k}$ , the terms in formula (3.7.2) that depend on  $p$  cancel out, one gets  $\Pr[A|x=k] = 1/\binom{n}{k}$ . Again there is no  $p$  in that formula.  $\square$

**PROBLEM 56.** *You perform a Bernoulli experiment, i.e., an experiment which can only have two outcomes, success  $s$  and failure  $f$ . The probability of success is  $p$ .*

• a. 3 points *You make 4 independent trials. Show that the probability that the first trial is successful, given that the total number of successes in the 4 trials is 3, is  $3/4$ .*

**ANSWER.** Let  $B = \{sfff, sffs, sfsf, sfss, ssff, ssfs, sssf, ssss\}$  be the event that the first trial is successful, and let  $\{x=3\} = \{fsss, sffs, ssfs, sssf\}$  be the event that there are 3 successes,

it has  $\binom{4}{3} = 4$  elements. Then

$$(3.7.3) \quad \Pr[B|x=3] = \frac{\Pr[B \cap \{x=3\}]}{\Pr[x=3]}$$

Now  $B \cap \{x=3\} = \{sffs, sffs, sssf\}$ , which has 3 elements. Therefore we get

$$(3.7.4) \quad \Pr[B|x=3] = \frac{3 \cdot p^3(1-p)}{4 \cdot p^3(1-p)} = \frac{3}{4}.$$

• b. 2 points *Discuss this result.*

**ANSWER.** It is significant that this probability is independent of  $p$ . I.e., once we know how many successes there were in the 4 trials, knowing the true  $p$  does not help us computing the probability of the event. From this also follows that the outcome of the event has no information about  $p$ . The value  $3/4$  is the same as the unconditional probability if  $p = 3/4$ . I.e., whether we know that the true frequency, the one that holds in the long run, is  $3/4$ , or whether we know that the actual frequency in this sample is  $3/4$ , both will lead us to the same predictions regarding the first throw. But not all conditional probabilities are equal to their unconditional counterparts: the conditional probability to get 3 successes in the first 4 trials is 1, but the unconditional probability is of course not 1.

### 3.8. Pitfalls of Data Reduction: The Ecological Fallacy

The nineteenth-century sociologist Emile Durkheim collected data on the frequency of suicides and the religious makeup of many contiguous provinces in Western Europe. He found that, on the average, provinces with greater proportions of Protestants had higher suicide rates and those with greater proportions of Catholics had lower suicide rates. Durkheim concluded from this that Protestants are more likely to commit suicide than Catholics. But this is not a compelling conclusion. It might have been that Catholics in predominantly Protestant provinces were taking their own lives. The oversight of this logical possibility is called the “Ecological Fallacy” [Sel58].

This seems like a far-fetched example, but arguments like this have been used to discredit data establishing connections between alcoholism and unemployment rates as long as the unit of investigation is not the individual but some aggregate.

One study [RZ78] found a positive correlation between driver education and the incidence of fatal automobile accidents involving teenagers. Closer analysis showed that the net effect of driver education was to put more teenagers on the road and therefore to increase rather than decrease the number of fatal crashes involving teenagers.

**PROBLEM 57.** 4 points *Assume your data show that counties with high rates of unemployment also have high rates of heart attacks. Can one conclude from this that the unemployed have a higher risk of heart attack? Discuss, besides the “ecological fallacy,” also other objections which one might make against such a conclusion.*

ANSWER. Ecological fallacy says that such a conclusion is only legitimate if one has individual data. Perhaps a rise in unemployment is associated with increased pressure and increased workloads among the employed, therefore it is the employed, not the unemployed, who get the heart attacks. Even if one has individual data one can still raise the following objection: perhaps unemployment and heart attacks are both consequences of a third variable (both unemployment and heart attacks depend on age or education, or freezing weather in a farming community causes unemployment for workers and heart attacks for the elderly).

□

But it is also possible to commit the opposite error and rely too much on individual data and not enough on “neighborhood effects.” In a relationship between health and income, it is much more detrimental for your health if you are poor in a poor neighborhood, than if you are poor in a rich neighborhood; and even wealthy people in a poor neighborhood do not escape some of the health and safety risks associated with this neighborhood.

Another pitfall of data reduction is Simpson’s paradox. According to table 1, the new drug was better than the standard drug both in urban and rural areas. But if you aggregate over urban and rural areas, then it looks like the standard drug was better than the new drug. This is an artificial example from [Spr98, p. 360].

Responses in Urban and Rural Areas to Each of Two Drugs				
	Standard Drug		New Drug	
	Urban	Rural	Urban	Rural
No Effect	500	350	1050	120
Cure	100	350	359	180

TABLE 1. Disaggregated Results of a New Drug

Response to Two Drugs		
	Standard Drug	New Drug
No Effect	850	1170
Cure	450	530

TABLE 2. Aggregated Version of Table 1

### 3.9. Independence of Random Variables

The concept of independence can be extended to random variables:  $x$  and  $y$  are independent if all events that can be defined in terms of  $x$  are independent of all events that can be defined in terms of  $y$ , i.e., all events of the form  $\{\omega \in U : x(\omega) \in C\}$  are independent of all events of the form  $\{\omega \in U : y(\omega) \in D\}$  with arbitrary (measurable)

subsets  $C, D \subset \mathbb{R}$ . Equivalent to this is that all events of the sort  $x \leq a$  are independent of all events of the sort  $y \leq b$ .

PROBLEM 58. 3 points *The simplest random variables are indicator functions, i.e., functions which can only take the values 0 and 1. Assume  $x$  is indicator function of the event  $A$  and  $y$  indicator function of the event  $B$ , i.e.,  $x$  takes the value 1 if  $A$  occurs, and the value 0 otherwise, and similarly with  $y$  and  $B$ . Show that according to the above definition of independence,  $x$  and  $y$  are independent if and only if events  $A$  and  $B$  are independent. (Hint: which are the only two events, other than the certain event  $U$  and the null event  $\emptyset$ , that can be defined in terms of  $x$ )?*

ANSWER. Only  $A$  and  $A'$ . Therefore we merely need the fact, shown in Problem 35, that if  $A$  and  $B$  are independent, then also  $A$  and  $B'$  are independent. By the same argument, also  $A'$  and  $B$  are independent, and  $A'$  and  $B'$  are independent. This is all one needs, except the observation that every event is independent of the certain event and the null event.

### 3.10. Location Parameters and Dispersion Parameters of a Random Variable

3.10.1. **Measures of Location.** A location parameter of random variables is a parameter which increases by  $c$  if one adds the constant  $c$  to the random variable.

The *expected value* is the most important location parameter. To motivate this, assume  $x$  is a discrete random variable, i.e., it takes the values  $x_1, \dots, x_r$  with probabilities  $p_1, \dots, p_r$  which sum up to one:  $\sum_{i=1}^r p_i = 1$ .  $x$  is observed  $n$  times independently. What can we expect the average value of  $x$  to be? For this we first need a formula for this average: if  $k_i$  is the number of times that  $x$  assumed the value  $x_i$  ( $i = 1, \dots, r$ ) then  $\sum k_i = n$ , and the average is  $\frac{k_1}{n}x_1 + \dots + \frac{k_r}{n}x_r$ . With appropriate definition of convergence, the relative frequencies  $\frac{k_i}{n}$  converge towards  $p_i$ . Therefore the average converges towards  $p_1x_1 + \dots + p_nx_n$ . This limit is the expected value of  $x$ , written as

$$(3.10.1) \quad E[x] = p_1x_1 + \dots + p_nx_n.$$

PROBLEM 59. *Why can one not use the usual concept of convergence here?*

ANSWER. Because there is no guarantee that the sample frequencies converge. It is not physically impossible (although it is highly unlikely) that certain outcome will never be realized.

Note the difference between the sample mean, i.e., the average measured in a given sample, and the “population mean” or expected value. The former is a random variable, the latter is a parameter. I.e., the former takes on a different value every time the experiment is performed, the latter does not.

Note that the expected value of the number of dots on a die is 3.5, which is not one of the possible outcomes when one rolls a die.

Expected value can be visualized as the center of gravity of the probability mass. If one of the tails has its weight so far out that there is no finite balancing point then the expected value is infinite or minus infinite. If both tails have their weights so far out that neither one has a finite balancing point, then the expected value does not exist.

It is trivial to show that for a function  $g(x)$  (which only needs to be defined for those values which  $x$  can assume with nonzero probability),  $E[g(x)] = p_1g(x_1) + \cdots + p_n g(x_n)$ .

Example of a countable probability mass distribution which has an infinite expected value:  $\Pr[x = x] = \frac{a}{x^2}$  for  $x = 1, 2, \dots$  ( $a$  is the constant  $1 / \sum_{i=1}^{\infty} \frac{1}{i^2}$ .) The expected value of  $x$  would be  $\sum_{i=1}^{\infty} \frac{a}{i}$ , which is infinite. But if the random variable is bounded, then its expected value exists.

The expected value of a *continuous* random variable is defined in terms of its density function:

$$(3.10.2) \quad E[x] = \int_{-\infty}^{+\infty} x f_x(x) dx$$

It can be shown that for any function  $g(x)$  defined for all those  $x$  for which  $f_x(x) \neq 0$  follows:

$$(3.10.3) \quad E[g(x)] = \int_{f_x(x) \neq 0} g(x) f_x(x) dx$$

Here the integral is taken over all the points which have nonzero density, instead of the whole line, because we did not require that the function  $g$  is defined at the points where the density is zero.

**PROBLEM 60.** *Let the random variable  $x$  have the Cauchy distribution, i.e., its density function is*

$$(3.10.4) \quad f_x(x) = \frac{1}{\pi(1+x^2)}$$

*Show that  $x$  does not have an expected value.*

ANSWER.

$$(3.10.5) \quad \int \frac{x dx}{\pi(1+x^2)} = \frac{1}{2\pi} \int \frac{2x dx}{1+x^2} = \frac{1}{2\pi} \int \frac{d(x^2)}{1+x^2} = \frac{1}{2\pi} \ln(1+x^2)$$

Rules about how to calculate with expected values (as long as they exist):

$$(3.10.6) \quad E[c] = c \text{ if } c \text{ is a constant}$$

$$(3.10.7) \quad E[ch] = cE[h]$$

$$(3.10.8) \quad E[h + j] = E[h] + E[j]$$

and if the random variables  $h$  and  $j$  are independent, then also

$$(3.10.9) \quad E[hj] = E[h]E[j].$$

**PROBLEM 61.** *2 points You make two independent trials of a Bernoulli experiment with success probability  $\theta$ , and you observe  $t$ , the number of successes. Compute the expected value of  $t^3$ . (Compare also Problem 169.)*

ANSWER.  $\Pr[t = 0] = (1 - \theta)^2$ ;  $\Pr[t = 1] = 2\theta(1 - \theta)$ ;  $\Pr[t = 2] = \theta^2$ . Therefore an application of (3.10.1) gives  $E[t^3] = 0^3 \cdot (1 - \theta)^2 + 1^3 \cdot 2\theta(1 - \theta) + 2^3 \cdot \theta^2 = 2\theta + 6\theta^2$ .

**THEOREM 3.10.1.** *Jensen's Inequality: Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be a function which is convex on an interval  $B \subset \mathbb{R}$ , which means*

$$(3.10.10) \quad g(\lambda a + (1 - \lambda)b) \leq \lambda g(a) + (1 - \lambda)g(b)$$

*for all  $a, b \in B$ . Furthermore let  $x : \mathbb{R} \rightarrow \mathbb{R}$  be a random variable so that  $\Pr[x \in B] = 1$ . Then  $g(E[x]) \leq E[g(x)]$ .*

**PROOF.** The Jensen inequality holds with equality if  $h(x)$  is a linear function (with a constant term), i.e., in this case,  $E[h(x)] = h(E[x])$ . (2) Therefore Jensen's inequality is proved if we can find a linear function  $h$  with the two properties  $h(E[x]) = g(E[x])$ , and  $h(x) \leq g(x)$  for all other  $x$ —because with such  $h$ ,  $E[g(x)] \geq E[h(x)] = h(E[x])$ . (3) The existence of such a  $h$  follows from convexity. Since  $g$  is convex, for every point  $a \in B$  there is a number  $\beta$  so that  $g(x) \geq g(a) + \beta(x - a)$ . This  $\beta$  is the slope of  $g$  if  $g$  is differentiable, and otherwise it is some number between the left and the right derivative (which both always exist for a convex function). We need this for  $a = E[x]$ .

This existence is the deepest part of this proof. We will not prove it here, for proof see [Rao73, pp. 57, 58]. One can view it as a special case of the separable hyperplane theorem.

**PROBLEM 62.** *Use Jensen's inequality to show that  $(E[x])^2 \leq E[x^2]$ . You are allowed to use, without proof, the fact that a function is convex on  $B$  if the second derivative exists on  $B$  and is nonnegative.*

**PROBLEM 63.** *Show that the expected value of the empirical distribution of a sample is the sample mean.*

□

Other measures of location: The *median* is that number  $m$  for which there is as much probability mass to the left of  $m$  as to the right, i.e.,

$$(3.10.11) \quad \Pr[x \leq m] = \frac{1}{2} \quad \text{or, equivalently,} \quad F_x(m) = \frac{1}{2}.$$

It is much more robust with respect to outliers than the mean. If there is more than one  $m$  satisfying (3.10.11), then some authors choose the smallest (in which case the median is a special case of the quantile function  $m = F^{-1}(1/2)$ ), and others the average between the biggest and smallest. If there is no  $m$  with property (3.10.11), i.e., if the cumulative distribution function jumps from a value that is less than  $\frac{1}{2}$  to a value that is greater than  $\frac{1}{2}$ , then the median is this jump point.

The *mode* is the point where the probability mass function or the probability density function is highest.

**3.10.2. Measures of Dispersion.** Here we will discuss *variance*, *standard deviation*, and *quantiles* and *percentiles*: The variance is defined as

$$(3.10.12) \quad \text{var}[x] = E[(x - E[x])^2],$$

but the formula

$$(3.10.13) \quad \text{var}[x] = E[x^2] - (E[x])^2$$

is usually more convenient.

How to calculate with variance?

$$(3.10.14) \quad \text{var}[ax] = a^2 \text{var}[x]$$

$$(3.10.15) \quad \text{var}[x + c] = \text{var}[x] \text{ if } c \text{ is a constant}$$

$$(3.10.16) \quad \text{var}[x + y] = \text{var}[x] + \text{var}[y] \text{ if } x \text{ and } y \text{ are independent.}$$

Note that the variance is additive only when  $x$  and  $y$  are independent; the expected value is always additive.

PROBLEM 64. Here we make the simple step from the definition of the variance to the usually more convenient formula (3.10.13).

• a. 2 points Derive the formula  $\text{var}[x] = E[x^2] - (E[x])^2$  from the definition of a variance, which is  $\text{var}[x] = E[(x - E[x])^2]$ . Hint: it is convenient to define  $\mu = E[x]$ . Write it down carefully, you will lose points for missing or unbalanced parentheses or brackets.

ANSWER. Here it is side by side with and without the notation  $E[x] = \mu$ :

$$(3.10.17) \quad \begin{aligned} \text{var}[x] &= E[(x - E[x])^2] & \text{var}[x] &= E[(x - \mu)^2] \\ &= E[x^2 - 2x(E[x]) + (E[x])^2] & &= E[x^2 - 2x\mu + \mu^2] \\ &= E[x^2] - 2(E[x])^2 + (E[x])^2 & &= E[x^2] - 2\mu^2 + \mu^2 \\ &= E[x^2] - (E[x])^2. & &= E[x^2] - \mu^2. \end{aligned}$$

• b. 1 point Assume  $\text{var}[x] = 3$ ,  $\text{var}[y] = 2$ ,  $x$  and  $y$  are independent. Compute  $\text{var}[-x]$ ,  $\text{var}[3y + 5]$ , and  $\text{var}[x - y]$ .

ANSWER. 3, 18, and 5.

PROBLEM 65. If all  $y_i$  are independent with same variance  $\sigma^2$ , then show that  $\bar{y}$  has variance  $\sigma^2/n$ .

The *standard deviation* is the square root of the variance. Often preferred because it has same scale as  $x$ . The variance, on the other hand, has the advantage of simple addition rule.

*Standardization*: if the random variable  $x$  has expected value  $\mu$  and standard deviation  $\sigma$ , then  $z = \frac{x - \mu}{\sigma}$  has expected value zero and variance one.

An  $\alpha$ th quantile or a 100 $\alpha$ th percentile of a random variable  $x$  was already defined previously to be the smallest number  $x$  so that  $\Pr[x \leq x] \geq \alpha$ .

**3.10.3. Mean-Variance Calculations.** If one knows mean and variance of a random variable, one does not by any means know the whole distribution, but one has already some information. For instance, one can compute  $E[y^2]$  from it, too.

PROBLEM 66. 4 points Consumer  $M$  has an expected utility function for money income  $u(x) = 12x - x^2$ . The meaning of an expected utility function is very simple: if he owns an asset that generates some random income  $y$ , then the utility he derives from this asset is the expected value  $E[u(y)]$ . He is contemplating acquiring two assets. One asset yields an income of 4 dollars with certainty. The other yields an expected income of 5 dollars with standard deviation 2 dollars. Does he prefer the certain or the uncertain asset?

ANSWER.  $E[u(y)] = 12E[y] - E[y^2] = 12E[y] - \text{var}[y] - (E[y])^2$ . Therefore the certain asset gives him utility  $48 - 0 - 16 = 32$ , and the uncertain one  $60 - 4 - 25 = 31$ . He prefers the certain asset.

**3.10.4. Moment Generating Function and Characteristic Function.** In this section we will use the exponential function  $e^x$ , also often written  $\exp(x)$ , which has the following properties:  $e^x = \lim_{n \rightarrow \infty} (1 + \frac{x}{n})^n$  (Euler's limit), and  $e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$

Many (but not all) random variables  $\mathbf{x}$  have a *moment generating function*  $m_{\mathbf{x}}(t)$  for certain values of  $t$ . If they do for  $t$  in an open interval around zero, then their distribution is uniquely determined by it. The definition is

$$(3.10.18) \quad m_{\mathbf{x}}(t) = \mathbb{E}[e^{t\mathbf{x}}]$$

It is a powerful computational device.

The moment generating function is in many cases a more convenient characterization of the random variable than the density function. It has the following uses:

1. One obtains the moments of  $\mathbf{x}$  by the simple formula

$$(3.10.19) \quad \mathbb{E}[\mathbf{x}^k] = \left. \frac{d^k}{dt^k} m_{\mathbf{x}}(t) \right|_{t=0}.$$

Proof:

$$(3.10.20) \quad e^{t\mathbf{x}} = 1 + t\mathbf{x} + \frac{t^2\mathbf{x}^2}{2!} + \frac{t^3\mathbf{x}^3}{3!} + \dots$$

$$(3.10.21) \quad m_{\mathbf{x}}(t) = \mathbb{E}[e^{t\mathbf{x}}] = 1 + t\mathbb{E}[\mathbf{x}] + \frac{t^2}{2!}\mathbb{E}[\mathbf{x}^2] + \frac{t^3}{3!}\mathbb{E}[\mathbf{x}^3] + \dots$$

$$(3.10.22) \quad \frac{d}{dt} m_{\mathbf{x}}(t) = \mathbb{E}[\mathbf{x}] + t\mathbb{E}[\mathbf{x}^2] + \frac{t^2}{2!}\mathbb{E}[\mathbf{x}^3] + \dots$$

$$(3.10.23) \quad \frac{d^2}{dt^2} m_{\mathbf{x}}(t) = \mathbb{E}[\mathbf{x}^2] + t\mathbb{E}[\mathbf{x}^3] + \dots \quad \text{etc.}$$

2. The moment generating function is also good for determining the probability distribution of linear combinations of independent random variables.

- a. it is easy to get the m.g.f. of  $\lambda\mathbf{x}$  from the one of  $\mathbf{x}$ :

$$(3.10.24) \quad m_{\lambda\mathbf{x}}(t) = m_{\mathbf{x}}(\lambda t)$$

because both sides are  $\mathbb{E}[e^{\lambda t\mathbf{x}}]$ .

- b. If  $\mathbf{x}$ ,  $\mathbf{y}$  independent, then

$$(3.10.25) \quad m_{\mathbf{x}+\mathbf{y}}(t) = m_{\mathbf{x}}(t)m_{\mathbf{y}}(t).$$

The proof is simple:

$$(3.10.26) \quad \mathbb{E}[e^{t(\mathbf{x}+\mathbf{y})}] = \mathbb{E}[e^{t\mathbf{x}}e^{t\mathbf{y}}] = \mathbb{E}[e^{t\mathbf{x}}]\mathbb{E}[e^{t\mathbf{y}}] \quad \text{due to independence.}$$

The *characteristic function* is defined as  $\psi_{\mathbf{x}}(t) = \mathbb{E}[e^{it\mathbf{x}}]$ , where  $i = \sqrt{-1}$ . It has the disadvantage that it involves complex numbers, but it has the advantage that it always exists, since  $\exp(ix) = \cos x + i \sin x$ . Since  $\cos$  and  $\sin$  are both bounded, they always have an expected value.

And, as its name says, the characteristic function characterizes the probability distribution. Analytically, many of its properties are similar to those of the moment generating function.

### 3.11. Entropy

**3.11.1. Definition of Information.** Entropy is the *average* information gain by the performance of the experiment. The *actual* information yielded by an event  $\mathbf{A}$  with probability  $\Pr[\mathbf{A}] = p \neq 0$  is defined as follows:

$$(3.11.1) \quad I[\mathbf{A}] = \log_2 \frac{1}{\Pr[\mathbf{A}]}$$

This is simply a transformation of the probability, and it has the dual interpretation of either how unexpected the event was, or the information yielded by the occurrence of event  $\mathbf{A}$ . It is characterized by the following properties [AD75, pp. 3–5]:

- $I[\mathbf{A}]$  only depends on the probability of  $\mathbf{A}$ , in other words, the information content of a message is independent of how the information is coded.
- $I[\mathbf{A}] \geq 0$  (nonnegativity), i.e., after knowing whether  $\mathbf{A}$  occurred we are more ignorant than before.
- If  $\mathbf{A}$  and  $\mathbf{B}$  are independent then  $I[\mathbf{A} \cap \mathbf{B}] = I[\mathbf{A}] + I[\mathbf{B}]$  (additivity independent events). This is the most important property.
- Finally the (inessential) normalization that if  $\Pr[\mathbf{A}] = 1/2$  then  $I[\mathbf{A}] = 1$ , i.e., a yes-or-no decision with equal probability (coin flip) is one unit of information.

Note that the information yielded by occurrence of the certain event is 0, and the information yielded by occurrence of the impossible event is  $\infty$ .

But the important information-theoretic results refer to average, not actual information, therefore let us define now *entropy*:

**3.11.2. Definition of Entropy.** The entropy of a probability field (experiment) is a measure of the uncertainty prevailing before the experiment is performed or of the *average* information yielded by the performance of this experiment. If the set  $\mathbf{U}$  of possible outcomes of the experiment has only a finite number of different elements, say their number is  $n$ , and the probabilities of these outcomes are  $p_1, \dots, p_n$ , then the Shannon entropy  $\mathbb{H}[\mathcal{F}]$  of this experiment is defined as

$$(3.11.2) \quad \frac{\mathbb{H}[\mathcal{F}]}{\text{bits}} = \sum_{k=1}^n p_k \log_2 \frac{1}{p_k}$$

This formula uses  $\log_2$ , logarithm with base 2, which can easily be computed from the natural logarithms,  $\log_2 x = \log x / \log 2$ . The choice of base 2 is convenient because in this way the most informative Bernoulli experiment, that with success probability  $p = 1/2$  (coin flip), has entropy 1. This is why one says: “the entropy is measured in bits.” If one goes over to logarithms of a different base, this simply means that one measures entropy in different units. In order to indicate this dependence on the measuring unit, equation (3.11.2) was written as the definition  $\frac{\mathbb{H}[\mathcal{F}]}{\text{bits}}$  instead of  $\mathbb{H}[\mathcal{F}]$ .



itself, i.e., this is the number one gets if one measures the entropy in bits. If one uses natural logarithms, then the entropy is measured in “nats.”

Entropy can be characterized axiomatically by the following axioms [K<sub>hi</sub>57]:

- The uncertainty associated with a finite complete scheme takes its largest value if all events are equally likely, i.e.,  $H(p_1, \dots, p_n) \leq H(1/n, \dots, 1/n)$ .
- The addition of an impossible event to a scheme does not change the amount of uncertainty.
- *Composition Law:* If the possible outcomes are arbitrarily combined into  $m$  groups  $W_1 = X_{11} \cup \dots \cup X_{1k_1}$ ,  $W_2 = X_{21} \cup \dots \cup X_{2k_2}$ ,  $\dots$ ,  $W_m = X_{m1} \cup \dots \cup X_{mk_m}$ , with corresponding probabilities  $w_1 = p_{11} + \dots + p_{1k_1}$ ,  $w_2 = p_{21} + \dots + p_{2k_2}$ ,  $\dots$ ,  $w_m = p_{m1} + \dots + p_{mk_m}$ , then

$$\begin{aligned}
H(p_1, \dots, p_n) &= H(w_1, \dots, w_m) + \\
&\quad + w_1 H(p_{11}/w_1 + \dots + p_{1k_1}/w_1) + \\
&\quad + w_2 H(p_{21}/w_2 + \dots + p_{2k_2}/w_2) + \dots + \\
&\quad + w_m H(p_{m1}/w_m + \dots + p_{mk_m}/w_m).
\end{aligned}$$

Since  $p_{ij}/w_j = \Pr[X_{ij}|W_j]$ , the composition law means: if you first learn half the outcome of the experiment, and then the other half, you will in the average get as much information as if you had been told the total outcome all at once.

The entropy of a *random variable*  $x$  is simply the entropy of the probability field induced by  $x$  on  $\mathbb{R}$ . It does not depend on the values  $x$  takes but only on the probabilities. For discretely distributed random variables it can be obtained by the following “eerily self-referential” prescription: plug the random variable into its own probability mass function and compute the expected value of the negative logarithm of this, i.e.,

$$(3.11.3) \quad \frac{H[x]}{\text{bits}} = E[-\log_2 p_x(x)]$$

One interpretation of the entropy is: it is the average number of yes-or-no questions necessary to describe the outcome of the experiment. For instance, consider an experiment which has 32 different outcomes occurring with equal probabilities. The entropy is

$$(3.11.4) \quad \frac{H}{\text{bits}} = \sum_{i=1}^{32} \frac{1}{32} \log_2 32 = \log_2 32 = 5 \quad \text{i.e.,} \quad H = 5 \text{ bits}$$

which agrees with the number of bits necessary to describe the outcome.

**PROBLEM 67.** Design a questioning scheme to find out the value of an integer between 1 and 32, and compute the expected number of questions in your scheme if all numbers are equally likely.

**ANSWER.** In binary digits one needs a number of length 5 to describe a number between 0 and 31, therefore the 5 questions might be: write down the binary expansion of your number minus zero, then: is the second binary digit in this expansion a zero, etc. Formulated without the use of binary digits these same questions would be: is the number between 1 and 16?, then: is it between 1 and 8 or 17 and 24?, then, is it between 1 and 4 or 9 or 12 or 17 and 20 or 25 and 28?, etc., the last question being whether it is odd. Of course, you can formulate those questions conditionally: First: between 1 and 16? if no, then second: between 1 and 8? if yes, then second: between 1 and 4? if yes, then second: between 1 and 2? if yes, then second: between 1 and 8? Etc. Each of these questions gives you exactly an entropy of 1 bit.

**PROBLEM 68.** [CT91, example 1.1.2 on p. 5] Assume there is a horse race with eight horses taking part. The probabilities for winning for the eight horses are  $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}$ .

- a. 1 point Show that the entropy of the horse race is 2 bits.

**ANSWER.**

$$\begin{aligned}
\frac{H}{\text{bits}} &= \frac{1}{2} \log_2 2 + \frac{1}{4} \log_2 4 + \frac{1}{8} \log_2 8 + \frac{1}{16} \log_2 16 + \frac{4}{64} \log_2 64 = \\
&= \frac{1}{2} + \frac{1}{2} + \frac{3}{8} + \frac{1}{4} + \frac{3}{8} = \frac{4+4+3+2+3}{8} = 2
\end{aligned}$$

- b. 1 point Suppose you want to send a binary message to another person indicating which horse won the race. One alternative is to assign the bit strings 000, 001, 010, 011, 100, 101, 110, 111 to the eight horses. This description requires 3 bits for any of the horses. But since the win probabilities are not uniform, it makes sense to use shorter descriptions for the horses more likely to win, so that we achieve a lower expected value of the description length. For instance, we could use the following set of bit strings for the eight horses: 0, 10, 110, 1110, 111100, 111110, 1111110, 1111111. Show that the the expected length of the message you send to your friend is 2 bits, as opposed to 3 bits for the uniform code. Note that in this case the expected value of the description length is equal to the entropy.

**ANSWER.** The math is the same as in the first part of the question:

$$\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{16} \cdot 4 + 4 \cdot \frac{1}{64} \cdot 6 = \frac{1}{2} + \frac{1}{2} + \frac{3}{8} + \frac{1}{4} + \frac{3}{8} = \frac{4+4+3+2+3}{8} = 2$$

**PROBLEM 69.** [CT91, example 2.1.2 on pp. 14/15]: The experiment has four possible outcomes; outcome  $x=a$  occurs with probability  $1/2$ ,  $x=b$  with probability  $1/4$ ,  $x=c$  with probability  $1/8$ , and  $x=d$  with probability  $1/8$ .

- a. 2 points The entropy of this experiment (in bits) is one of the following three numbers:  $11/8$ ,  $7/4$ ,  $2$ . Which is it?

• b. 2 points Suppose we wish to determine the outcome of this experiment with the minimum number of questions. An efficient first question is “Is  $x=a$ ?” This splits the probability in half. If the answer to the first question is no, then the second question can be “Is  $x=b$ ?” The third question, if it is necessary, can then be: “Is  $x=c$ ?” Compute the expected number of binary questions required.

• c. 2 points Show that the entropy gained by each question is 1 bit.

• d. 3 points Assume we know about the first outcome that  $x \neq a$ . What is the entropy of the remaining experiment (i.e., under the conditional probability)?

• e. 5 points Show in this example that the composition law for entropy holds.

PROBLEM 70. 2 points In terms of natural logarithms equation (3.11.4) defining entropy reads

$$(3.11.5) \quad \frac{H}{\text{bits}} = \frac{1}{\ln 2} \sum_{k=1}^n p_k \ln \frac{1}{p_k}.$$

Compute the entropy of (i.e., the average information gained by) a roll of an unbiased die.

ANSWER. Same as the actual information gained, since each outcome is equally likely:

$$(3.11.6) \quad \frac{H}{\text{bits}} = \frac{1}{\ln 2} \left( \frac{1}{6} \ln 6 + \cdots + \frac{1}{6} \ln 6 \right) = \frac{\ln 6}{\ln 2} = 2.585$$

□

• a. 3 points How many questions does one need in the average to determine the outcome of the roll of an unbiased die? In other words, pick a certain questioning scheme (try to make it efficient) and compute the average number of questions if this scheme is followed. Note that this average cannot be smaller than the entropy  $H/\text{bits}$ , and if one chooses the questions optimally, it is smaller than  $H/\text{bits} + 1$ .

ANSWER. First question: is it bigger than 3? Second question: is it even? Third question (if necessary): is it a multiple of 3? In this scheme, the number of questions for the six faces of the die are 3, 2, 3, 3, 2, 3, therefore the average is  $\frac{4}{6} \cdot 3 + \frac{2}{6} \cdot 2 = 2\frac{2}{3}$ . Also optimal: (1) is it bigger than 2? (2) is it odd? (3) is it bigger than 4? Gives 2, 2, 3, 3, 3, 3. Also optimal: 1st question: is it 1 or 2? If answer is no, then second question is: is it 3 or 4?; otherwise go directly to the third question: is it odd or even? The steamroller approach: Is it 1? Is it 2? etc. gives 1, 2, 3, 4, 5, 5 with expected number  $3\frac{1}{3}$ . Even this is here  $< 1 + H/\text{bits}$ . □

PROBLEM 71.

• a. 1 point Compute the entropy of a roll of two unbiased dice if they are distinguishable.

ANSWER. Just twice the entropy from Problem 70.

$$(3.11.7) \quad \frac{H}{\text{bits}} = \frac{1}{\ln 2} \left( \frac{1}{36} \ln 36 + \cdots + \frac{1}{36} \ln 36 \right) = \frac{\ln 36}{\ln 2} = 5.170$$

• b. Would you expect the entropy to be greater or less in the more usual case that the dice are indistinguishable? Check your answer by computing it.

ANSWER. If the dice are indistinguishable, then one gets less information, therefore the experiment has less entropy. One has six like pairs with probability  $1/36$  and  $6 \cdot 5/2 = 15$  unlike pairs with probability  $2/36 = 1/18$  each. Therefore the average information gained is

$$(3.11.8) \quad \frac{H}{\text{bits}} = \frac{1}{\ln 2} \left( 6 \cdot \frac{1}{36} \ln 36 + 15 \cdot \frac{1}{18} \ln 18 \right) = \frac{1}{\ln 2} \left( \frac{1}{6} \ln 36 + \frac{5}{6} \ln 18 \right) = 4.337$$

• c. 3 points Note that the difference between these two entropies is  $5/6 = 0.833$ . How can this be explained?

ANSWER. This is the composition law (??) in action. Assume you roll two dice which you first consider indistinguishable and afterwards someone tells you which is which. How much information do you gain? Well, if the numbers are the same, then telling you which die is which does not give you any information, since the outcomes of the experiment are defined as: which number has the first die, which number has the second die, regardless of where on the table the dice land. But if the numbers are different, then telling you which is which allows you to discriminate between two outcomes both of which have conditional probability  $1/2$  given the outcome you already know. In this case the information you gain is therefore 1 bit. Since the probability of getting two different numbers is  $5/6$ , the expected value of the information gained explains the difference in entropy.

All these definitions use the convention  $0 \log \frac{1}{0} = 0$ , which can be justified by the following continuity argument: Define the function, graphed in Figure 3:

$$(3.11.9) \quad \eta(w) = \begin{cases} w \log \frac{1}{w} & \text{if } w > 0 \\ 0 & \text{if } w = 0. \end{cases}$$

$\eta$  is continuous for all  $w \geq 0$ , even at the boundary point  $w = 0$ . Differentiation gives  $\eta'(w) = -(1 + \log w)$ , and  $\eta''(w) = -w^{-1}$ . The function starts out at the origin with a vertical tangent, and since the second derivative is negative, it is strictly concave for all  $w > 0$ . The definition of strict concavity is  $\eta(w) < \eta(v) + (w - v)\eta'(v)$  for  $w \neq v$ , i.e., the function lies below all its tangents. Substituting  $\eta'(v) = -(1 + \log v)$  and simplifying gives  $w - w \log w \leq v - w \log v$  for  $v, w > 0$ . One verifies that this inequality also holds for  $v, w \geq 0$ .

PROBLEM 72. Make a complete proof, discussing all possible cases, that  $w - w \log w \leq v - w \log v$  follows

$$(3.11.10) \quad w - w \log w \leq v - w \log v$$

ANSWER. We already know it for  $v, w > 0$ . Now if  $v = 0$  and  $w = 0$  then the equation reads  $0 \leq 0$ ; if  $v > 0$  and  $w = 0$  the equation reads  $0 \leq v$ , and if  $w > 0$  and  $v = 0$  then the equation reads  $w - w \log w \leq +\infty$ .  $\square$

**3.11.3. How to Keep Forecasters Honest.** This mathematical result allows an interesting alternative mathematical characterization of entropy. Assume Anita performs a Bernoulli experiment whose success probability she does not know but wants to know. Clarence knows this probability but is not on very good terms with Anita; therefore Anita is unsure that he will tell the truth if she asks him.

Anita knows “how to keep forecasters honest.” She proposes the following deal to Clarence: “you tell me the probability  $q$ , and after performing my experiment I pay you the amount  $\log_2(q)$  if the experiment is a success, and  $\log_2(1 - q)$  if it is a failure. If Clarence agrees to this deal, then telling Anita that value  $q$  which is the true success probability of the Bernoulli experiment maximizes the expected value of his payoff. And the maximum expected value of this payoff is exactly the negative of the entropy of the experiment.

Proof: Assume the correct value of the probability is  $p$ , and the number Clarence tells Tina is  $q$ . For every  $p, q$  between 0 and 1 we have to show:

$$(3.11.11) \quad p \log p + (1 - p) \log(1 - p) \geq p \log q + (1 - p) \log(1 - q).$$

For this, plug  $w = p$  and  $v = q$  as well as  $w = 1 - p$  and  $v = 1 - q$  into equation (3.11.10) and add.

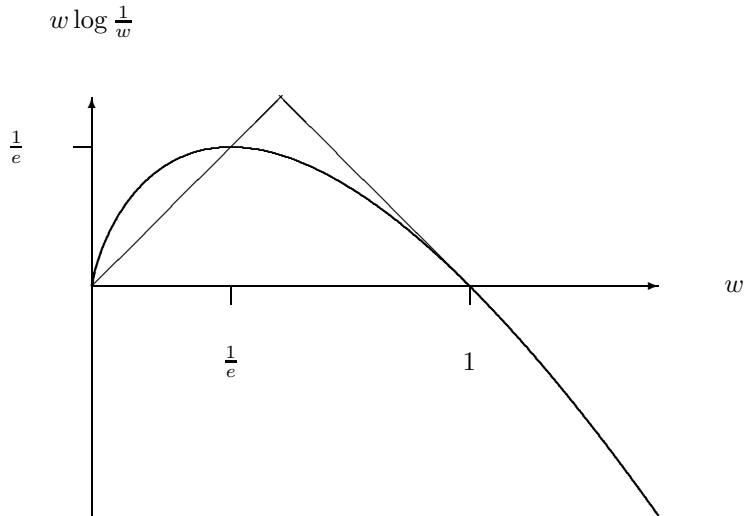


FIGURE 3.  $\eta : w \mapsto w \log \frac{1}{w}$  is continuous at 0, and concave everywhere

**3.11.4. The Inverse Problem.** Now let us go over to the inverse problem: computing those probability fields which have maximum entropy subject to the information you have.

If you know that the experiment has  $n$  different outcomes, and you do not know the probabilities of these outcomes, then the maximum entropy approach amounts to assigning equal probability  $1/n$  to each outcome.

PROBLEM 73. (Not eligible for in-class exams) You are playing a slot machine. Feeding one dollar to this machine leads to one of four different outcomes:  $E_1$ : machine returns nothing, i.e., you lose \$1.  $E_2$ : machine returns \$1, i.e., you lose nothing and win nothing.  $E_3$ : machine returns \$2, i.e., you win \$1.  $E_4$ : machine returns \$9, i.e., you win \$9. Events  $E_i$  occurs with probability  $p_i$ , but these probabilities are unknown. But due to a new “Truth-in-Gambling Act” you find a sticker on the side of the machine which says that in the long run the machine pays out only \$0.90 for every dollar put in. Show that those values of  $p_1, p_2, p_3$ , and  $p_4$  which maximize the entropy (and therefore make the machine most interesting) subject to the constraint that the expected payoff per dollar put in is \$0.90, are  $p_1 = 0.4473$ ,  $p_2 = 0.3118$ ,  $p_3 = 0.2231$ ,  $p_4 = 0.0138$ .

ANSWER. Solution is derived in [Rie85, pp. 68/9 and 74/5], and he refers to [Rie77]. You have to maximize  $-\sum p_n \log p_n$  subject to  $\sum p_n = 1$  and  $\sum c_n p_n = d$ . In our case  $c_1 = 0$ ,  $c_2 = 0$ ,  $c_3 = 2$ , and  $c_4 = 10$ , and  $d = 0.9$ , but the treatment below goes through for arbitrary  $c_i$  as long as not all of them are equal. This case is discussed in detail in the answer to Problem 74.

• a. *Difficult:* Does the maximum entropy approach also give us some guidelines how to select these probabilities if all we know is that the expected value of the payoff rate is smaller than 1?

ANSWER. As shown in [Rie85, pp. 68/9 and 74/5], one can give the minimum value of entropy for all distributions with payoff smaller than 1:  $H < 1.6590$ , and one can also give some bounds for the probabilities:  $p_1 > 0.4272$ ,  $p_2 < 0.3167$ ,  $p_3 < 0.2347$ ,  $p_4 < 0.0214$ .

• b. What if you also know that the entropy of this experiment is 1.5?

ANSWER. This was the purpose of the paper [Rie85].

PROBLEM 74. (Not eligible for in-class exams) Let  $p_1, p_2, \dots, p_n$  ( $\sum p_i = 1$ ) be the proportions of the population of a city living in  $n$  residential colonies. The cost of living in colony  $i$ , which includes cost of travel from the colony to the central business district, the cost of the time this travel consumes, the rent or mortgage payments and other costs associated with living in colony  $i$ , is represented by the monetary amount  $c_i$ . Without loss of generality we will assume that the  $c_i$  are numbered such a way that  $c_1 \leq c_2 \leq \dots \leq c_n$ . We will also assume that the  $c_i$  are not equal. We assume that the  $c_i$  are known and that also the average expenditures for travel etc. in the population is known; its value is  $d$ . One approach to modelling

population distribution is to maximize the entropy subject to the average expenditures, i.e., to choose  $p_1, p_2, \dots, p_n$  such that  $H = \sum p_i \log \frac{1}{p_i}$  is maximized subject to the two constraints  $\sum p_i = 1$  and  $\sum p_i c_i = d$ . This would give the greatest uncertainty about where someone lives.

- a. 3 points Set up the Lagrange function and show that

$$(3.11.12) \quad p_i = \frac{\exp(-\lambda c_i)}{\sum \exp(-\lambda c_i)}$$

where the Lagrange multiplier  $\lambda$  must be chosen such that  $\sum p_i c_i = d$ .

ANSWER. The Lagrange function is

$$(3.11.13) \quad L = - \sum p_n \log p_n - \kappa (\sum p_n - 1) - \lambda (\sum c_n p_n - d)$$

Partial differentiation with respect to  $p_i$  gives the first order conditions

$$(3.11.14) \quad -\log p_i - 1 - \kappa - \lambda c_i = 0.$$

Therefore  $p_i = \exp(-\kappa - 1) \exp(-\lambda c_i)$ . Plugging this into the first constraint gives  $1 = \sum p_i = \exp(-\kappa - 1) \sum \exp(-\lambda c_i)$  or  $\exp(-\kappa - 1) = \frac{1}{\sum \exp(-\lambda c_i)}$ . This constraint therefore defines  $\kappa$  uniquely, and we can eliminate  $\kappa$  from the formula for  $p_i$ :

$$(3.11.15) \quad p_i = \frac{\exp(-\lambda c_i)}{\sum \exp(-\lambda c_i)}$$

Now all the  $p_i$  depend on the same unknown  $\lambda$ , and this  $\lambda$  must be chosen such that the second constraint holds. This is the Maxwell-Boltzmann distribution if  $\mu = kT$  where  $k$  is the Boltzmann constant and  $T$  the temperature.  $\square$

- b. 2 points Here is a mathematical lemma needed for the next part: Prove that for  $a_i \geq 0$  and  $c_i$  arbitrary follows  $\sum a_i \sum a_i c_i^2 \geq (\sum a_i c_i)^2$ , and if all  $a_i > 0$  and not all  $c_i$  equal, then this inequality is strict.

ANSWER. By choosing the same subscripts in the second sum as in the first we pair elements of the first sum with elements of the second sum:

$$(3.11.16) \quad \sum_i a_i \sum_j c_j^2 a_j - \sum_i c_i a_i \sum_j c_j a_j = \sum_{i,j} (c_j^2 - c_i c_j) a_i a_j$$

but if we interchange  $i$  and  $j$  on the rhs we get

$$(3.11.17) \quad = \sum_{j,i} (c_i^2 - c_j c_i) a_j a_i = \sum_{i,j} (c_i^2 - c_i c_j) a_i a_j$$

Now add the righthand sides to get

$$(3.11.18) \quad 2 \left( \sum_i a_i \sum_j c_j^2 a_j - \sum_i c_i a_i \sum_j c_j a_j \right) = \sum_{i,j} (c_i^2 + c_j^2 - 2c_i c_j) a_i a_j = \sum_{i,j} (c_i - c_j)^2 a_i a_j \geq 0$$

$\square$

- c. 3 points It is not possible to solve equations (3.11.12) analytically for  $\lambda$ , the following can be shown [Kap89, p. 310/11]: the function  $f$  defined by

$$(3.11.19) \quad f(\lambda) = \frac{\sum c_i \exp(-\lambda c_i)}{\sum \exp(-\lambda c_i)}$$

is a strictly decreasing function which decreases from  $c_n$  to  $c_1$  as  $\lambda$  goes from  $-\infty$  to  $\infty$ , and  $f(0) = \bar{c}$  where  $\bar{c} = (1/n) \sum c_i$ . We need that  $\lambda$  for which  $f(\lambda) = d$ , and this equation has no real root if  $d < c_1$  or  $d > c_n$ , it has a unique positive root if  $c_1 < d < \bar{c}$  it has the unique root 0 for  $d = \bar{c}$ , and it has a unique negative root if  $\bar{c} < d < c_n$ . From this follows: as long as  $d$  lies between the lowest and highest cost and as long as the cost numbers are not all equal, the  $p_i$  are uniquely determined by the above entropy maximization problem.

ANSWER. Here is the derivative; it is negative because of the mathematical lemma just shown

$$(3.11.20) \quad f'(\lambda) = \frac{u'v - uv'}{v^2} = - \frac{\sum \exp(-\lambda c_i) \sum c_i^2 \exp(-\lambda c_i) - \left( \sum c_i \exp(-\lambda c_i) \right)^2}{\left( \sum \exp(-\lambda c_i) \right)^2} < 0$$

Since  $c_1 \leq c_2 \leq \dots \leq c_n$ , it follows

$$(3.11.21) \quad c_1 = \frac{\sum c_1 \exp(-\lambda c_i)}{\sum \exp(-\lambda c_i)} \leq \frac{\sum c_i \exp(-\lambda c_i)}{\sum \exp(-\lambda c_i)} \leq \frac{\sum c_n \exp(-\lambda c_i)}{\sum \exp(-\lambda c_i)} = c_n$$

Now the statement about the limit can be shown if not all  $c_j$  are equal, say  $c_1 < c_{k+1}$  but  $c_1 = c_2 = \dots = c_k$ . The fraction can be written as

$$(3.11.22) \quad \frac{k c_1 \exp(-\lambda c_1) + \sum_{i=1}^{n-k} c_{k+i} \exp(-\lambda c_{k+i})}{k \exp(-\lambda c_1) + \sum_{i=1}^{n-k} \exp(-\lambda c_{k+i})} = \frac{k c_1 + \sum_{i=1}^{n-k} c_{k+i} \exp(-\lambda(c_{k+i} - c_1))}{k + \sum_{i=1}^{n-k} \exp(-\lambda(c_{k+i} - c_1))}$$

Since  $c_{k+i} - c_1 > 0$ , this converges towards  $c_1$  for  $\lambda \rightarrow \infty$ .

- d. 3 points Show that the maximum attained entropy is  $H = \lambda d + k(\lambda)$  where

$$(3.11.23) \quad k(\lambda) = \log \left( \sum \exp(-\lambda c_j) \right).$$

Although  $\lambda$  depends on  $d$ , show that  $\frac{\partial H}{\partial d} = \lambda$ , i.e., it is the same as if  $\lambda$  did not depend on  $d$ . This is an example of the “envelope theorem,” and it also gives an interpretation of  $\lambda$ .

ANSWER. We have to plug the optimal  $p_i = \frac{\exp(-\lambda c_i)}{\sum \exp(-\lambda c_i)}$  into the formula for  $H = - \sum p_i \log p_i$

For this note that  $-\log p_i = \lambda c_i + k(\lambda)$  where  $k(\lambda) = \log(\sum \exp(-\lambda c_j))$  does not depend on  $d$ . Therefore  $H = \sum p_i (\lambda c_i + k(\lambda)) = \lambda \sum p_i c_i + k(\lambda) \sum p_i = \lambda d + k(\lambda)$ , and  $\frac{\partial H}{\partial d} = \lambda + d \frac{\partial \lambda}{\partial d} + k'(\lambda)$ . Now we need the derivative of  $k(\lambda)$ , and we discover that  $k'(\lambda) = -f(\lambda)$  where  $f(\lambda)$  was defined in (3.11.19). Therefore  $\frac{\partial H}{\partial d} = \lambda + (d - f(\lambda)) \frac{\partial \lambda}{\partial d} = \lambda$ .

- e. 5 points Now assume  $d$  is not known (but the  $c_i$  are still known), i.e., know that (3.11.12) holds for some  $\lambda$  but we don't know which. We want to estim

this  $\lambda$  (and therefore all  $p_i$ ) by taking a random sample of  $m$  people from that metropolitan area and asking them what their regional living expenditures are and where they live. Assume  $x_i$  people in this sample live in colony  $i$ . One way to estimate this  $\lambda$  would be to use the average consumption expenditure of the sample,  $\sum \frac{x_i}{m} c_i$ , as an estimate of the missing  $d$  in the above procedure, i.e., choose that  $\lambda$  which satisfies  $f(\lambda) = \sum \frac{x_i}{m} c_i$ . Another procedure, which seems to make a better use of the information given by the sample, would be to compute the maximum likelihood estimator of  $\lambda$  based on all  $x_i$ . Show that these two estimation procedures are identical.

ANSWER. The  $x_i$  have the multinomial distribution. Therefore, given that the proportion  $p_i$  of the population lives in colony  $i$ , and you are talking a random sample of size  $m$  from the whole population, then the probability to get the outcome  $x_1, \dots, x_n$  is

$$(3.11.24) \quad L = \frac{m!}{x_1! \cdots x_n!} p_1^{x_1} p_2^{x_2} \cdots p_n^{x_n}$$

This is what we have to maximize, subject to the condition that the  $p_i$  are an entropy maximizing population distribution. Let's take logs for computational simplicity:

$$(3.11.25) \quad \log L = \log m! - \sum_j \log x_j! + \sum x_i \log p_i$$

All we know about the  $p_i$  is that they must be some entropy maximizing probabilities, but we don't know yet which ones, i.e., they depend on the unknown  $\lambda$ . Therefore we need the formula again  $-\log p_i = \lambda c_i + k(\lambda)$  where  $k(\lambda) = \log(\sum \exp(-\lambda c_j))$  does not depend on  $i$ . This gives

$$(3.11.26) \quad \log L = \log m! - \sum_j \log x_j! - \sum x_i (\lambda c_i + k(\lambda)) = \log m! - \sum_j \log x_j! - \lambda \sum x_i c_i + k(\lambda) m$$

(for this last term remember that  $\sum x_i = m$ . Therefore the derivative is

$$(3.11.27) \quad \frac{1}{m} \frac{\partial}{\partial \lambda} \log L = \sum \frac{x_i}{m} c_i - f(\lambda)$$

I.e., using the obvious estimate for  $d$  is the same as maximum likelihood under the assumption of maximum entropy.  $\square$

This is a powerful estimation strategy. An article with sensational image re-constitutions using maximum entropy algorithms is [SG85, pp. 111, 112, 115, 116]. And [GJM96] applies maximum entropy methods to ill-posed or underdetermined problems in econometrics!

## CHAPTER 4

## Specific Random Variables

## 4.1. Binomial

We will begin with mean and variance of the binomial variable, i.e., the number of successes in  $n$  independent repetitions of a Bernoulli trial (3.7.1). The binomial variable has the two parameters  $n$  and  $p$ . Let us look first at the case  $n = 1$ , in which the binomial variable is also called *indicator variable*: If the event  $A$  has probability  $p$ , then its complement  $A'$  has the probability  $q = 1 - p$ . The indicator variable of  $A$ , which assumes the value 1 if  $A$  occurs, and 0 if it doesn't, has expected value  $p$  and variance  $pq$ . For the binomial variable with  $n$  observations, which is the sum of  $n$  independent indicator variables, the expected value (mean) is  $np$  and the variance is  $npq$ .

**PROBLEM 75.** *The random variable  $x$  assumes the value  $a$  with probability  $p$  and the value  $b$  with probability  $q = 1 - p$ . Show that  $\text{var}[x] = pq(a - b)^2$ .*

**ANSWER.**  $E[x] = pa + qb$ ;  $\text{var}[x] = E[x^2] - (E[x])^2 = pa^2 + qb^2 - (pa + qb)^2 = (p - p^2)a^2 - 2pqab + (q - q^2)b^2 = pq(a - b)^2$ . For this last equality we need  $p - p^2 = p(1 - p) = pq$ .  $\square$

The *Negative Binomial Variable* is, like the binomial variable, derived from the Bernoulli experiment; but one reverses the question. Instead of asking how many successes one gets in a given number of trials, one asks, how many trials one must make to get a given number of successes, say,  $r$  successes.

First look at  $r = 1$ . Let  $t$  denote the number of the trial at which the first success occurs. Then

$$(4.1.1) \quad \Pr[t=n] = pq^{n-1} \quad (n = 1, 2, \dots).$$

This is called the geometric probability.

Is the probability derived in this way  $\sigma$ -additive? The sum of a geometrically declining sequence is easily computed:

$$(4.1.2) \quad 1 + q + q^2 + q^3 + \dots = s \quad \text{Now multiply by } q:$$

$$(4.1.3) \quad q + q^2 + q^3 + \dots = qs \quad \text{Now subtract and write } 1 - q = p:$$

$$(4.1.4) \quad 1 = ps$$

Equation (4.1.4) means  $1 = p + pq + pq^2 + \dots$ , i.e., the sum of all probabilities indeed 1.

Now what is the expected value of a geometric variable? Use definition of expected value of a discrete variable:  $E[t] = p \sum_{k=1}^{\infty} kq^{k-1}$ . To evaluate the infinite sum, solve (4.1.4) for  $s$ :

$$(4.1.5) \quad s = \frac{1}{p} \quad \text{or} \quad 1 + q + q^2 + q^3 + q^4 \dots = \sum_{k=0}^{\infty} q^k = \frac{1}{1 - q}$$

and differentiate both sides with respect to  $q$ :

$$(4.1.6) \quad 1 + 2q + 3q^2 + 4q^3 + \dots = \sum_{k=1}^{\infty} kq^{k-1} = \frac{1}{(1 - q)^2} = \frac{1}{p^2}.$$

The expected value of the geometric variable is therefore  $E[t] = \frac{p}{p^2} = \frac{1}{p}$ .

**PROBLEM 76.** *Assume  $t$  is a geometric random variable with parameter  $p$ , i.e., it has the values  $k = 1, 2, \dots$  with probabilities*

$$(4.1.7) \quad p_t(k) = pq^{k-1}, \quad \text{where } q = 1 - p.$$

*The geometric variable denotes the number of times one has to perform a Bernoulli experiment with success probability  $p$  to get the first success.*

• a. *1 point Given a positive integer  $n$ . What is  $\Pr[t > n]$ ? (Easy with a simple trick!)*

**ANSWER.**  $t > n$  means, the first  $n$  trials must result in failures, i.e.,  $\Pr[t > n] = q^n$ . Since  $\{t > n\} = \{t = n + 1\} \cup \{t = n + 2\} \cup \dots$ , one can also get the same result in a more tedious way. It is  $pq^n + pq^{n+1} + pq^{n+2} + \dots = s$ , say. Therefore  $qs = pq^{n+1} + pq^{n+2} + \dots$ , and  $(1 - q)s = pq^n$  since  $p = 1 - q$ , it follows  $s = q^n$ .

• b. *2 points Let  $m$  and  $n$  be two positive integers with  $m < n$ . Show that  $\Pr[t=n|t>m] = \Pr[t=n - m]$ .*

$$\text{ANSWER. } \Pr[t=n|t>m] = \frac{\Pr[t=n]}{\Pr[t>m]} = \frac{pq^{n-1}}{q^m} = pq^{n-m-1} = \Pr[t=n - m].$$

• c. *1 point Why is this property called the memory-less property of the geometric random variable?*

**ANSWER.** If you have already waited for  $m$  periods without success, the probability that success will come in the  $n$ th period is the same as the probability that it comes in  $n - m$  periods if you start now. Obvious if you remember that geometric random variable is time you have to wait until 1st success in Bernoulli trial.

PROBLEM 77.  $t$  is a geometric random variable as in the preceding problem. In order to compute  $\text{var}[t]$  it is most convenient to make a detour via  $E[t(t-1)]$ . Here are the steps:

- a. Express  $E[t(t-1)]$  as an infinite sum.

ANSWER. Just write it down according to the definition of expected values:  $\sum_{k=0}^{\infty} k(k-1)pq^{k-1} = \sum_{k=2}^{\infty} k(k-1)pq^{k-1}$ .  $\square$

- b. Derive the formula

$$(4.1.8) \quad \sum_{k=2}^{\infty} k(k-1)q^{k-2} = \frac{2}{(1-q)^3}$$

by the same trick by which we derived a similar formula in class. Note that the sum starts at  $k=2$ .

ANSWER. This is just a second time differentiating the geometric series, i.e., first time differentiating (4.1.6).  $\square$

- c. Use a. and b. to derive

$$(4.1.9) \quad E[t(t-1)] = \frac{2q}{p^2}$$

ANSWER.

$$(4.1.10) \quad \sum_{k=2}^{\infty} k(k-1)pq^{k-1} = pq \sum_{k=2}^{\infty} k(k-1)q^{k-2} = pq \frac{2}{(1-q)^3} = \frac{2q}{p^2}.$$

- d. Use c. and the fact that  $E[t] = 1/p$  to derive

$$(4.1.11) \quad \text{var}[t] = \frac{q}{p^2}.$$

ANSWER.

$$(4.1.12) \quad \text{var}[t] = E[t^2] - (E[t])^2 = E[t(t-1)] + E[t] - (E[t])^2 = \frac{2q}{p^2} + \frac{1}{p} - \frac{1}{p^2} = \frac{q}{p^2}.$$

Now let us look at the negative binomial with arbitrary  $r$ . What is the probability that it takes  $n$  trials to get  $r$  successes? (That means, with  $n-1$  trials we did not yet have  $r$  successes.) The probability that the  $n$ th trial is a success is  $p$ . The probability that there are  $r-1$  successes in the first  $n-1$  trials is  $\binom{n-1}{r-1}p^{r-1}q^{n-r}$ . Multiply those to get:

$$(4.1.13) \quad \Pr[t=n] = \binom{n-1}{r-1}p^r q^{n-r}.$$

This is the negative binomial, also called the Pascal probability distribution with parameters  $r$  and  $p$ .

One easily gets the mean and variance, because due to the memory-less property it is the sum of  $r$  independent geometric variables:

$$(4.1.14) \quad E[t] = \frac{r}{p} \quad \text{var}[t] = \frac{rq}{p^2}$$

Some authors define the negative binomial as the number of failures before the  $r$ th success. Their formulas will look slightly different than ours.

PROBLEM 78. 3 points A fair coin is flipped until heads appear 10 times, and  $t$  is the number of times tails appear before the 10th appearance of heads. Show that the expected value  $E[x] = 10$ .

ANSWER. Let  $t$  be the number of the throw which gives the 10th head.  $t$  is a negative binomial with  $r=10$  and  $p=1/2$ , therefore  $E[t] = 20$ . Since  $x = t - 10$ , it follows  $E[x] = 10$ .

PROBLEM 79. (Banach's match-box problem) (Not eligible for in-class exam) There are two restaurants in town serving hamburgers. In the morning each of them obtains a shipment of  $n$  raw hamburgers. Every time someone in that town wants to eat a hamburger, he or she selects one of the two restaurants at random. What is the probability that the  $(n+k)$ th customer will have to be turned away because the restaurant selected has run out of hamburgers?

ANSWER. For each restaurant it is the negative binomial probability distribution in disguise. If a restaurant runs out of hamburgers this is like having  $n$  successes in  $n+k$  tries.

But one can also reason it out: Assume one of the restaurants must turn customers away after the  $n+k$ th customer. Write down all the  $n+k$  decisions made: write a 1 if the customer goes to the first restaurant, and a 2 if he goes to the second. I.e., write down  $n+k$  ones and twos. Under what conditions will such a sequence result in the  $n+k$ th move eating the last hamburger at the first restaurant? Exactly if it has  $n$  ones and  $k$  twos, a  $n+k$ th move is a one. As in the reasoning for the negative binomial probability distribution, there are  $\binom{n+k-1}{n-1}$  possibilities, each of which has probability  $2^{-(n+k)}$ . Emptying the second restaurant has the same probability. Together the probability is therefore  $\binom{n+k-1}{n-1}2^{1-n-k}$ .

## 4.2. The Hypergeometric Probability Distribution

Until now we had independent events, such as, repeated throwing of coins or dice, sampling with replacement from finite populations, or sampling from infinite populations. If we sample *without* replacement from a *finite* population, the probability of the second element of the sample depends on what the first element was. Here the hypergeometric probability distribution applies.

Assume we have an urn with  $w$  white and  $n-w$  black balls in it, and we take a sample of  $m$  balls. What is the probability that  $y$  of them are white?

We are not interested in the order in which these balls are taken out; we may therefore assume that they are taken out simultaneously, therefore the set  $\mathbf{U}$  of outcomes is the set of subsets containing  $m$  of the  $n$  balls. The total number of such subsets is  $\binom{n}{m}$ . How many of them have  $y$  white balls in them? Imagine you first pick  $y$  white balls from the set of all white balls (there are  $\binom{w}{y}$  possibilities to do that), and then you pick  $m - y$  black balls from the set of all black balls, which can be done in  $\binom{n-w}{m-y}$  different ways. Every union of such a set of white balls with a set of black balls gives a set of  $m$  elements with exactly  $y$  white balls, as desired. There are therefore  $\binom{w}{y}\binom{n-w}{m-y}$  different such sets, and the probability of picking such a set is

$$(4.2.1) \quad \Pr[\text{Sample of } m \text{ elements has exactly } y \text{ white balls}] = \frac{\binom{w}{y}\binom{n-w}{m-y}}{\binom{n}{m}}.$$

**PROBLEM 80.** *You have an urn with  $w$  white and  $n - w$  black balls in it, and you take a sample of  $m$  balls with replacement, i.e., after pulling each ball out you put it back in before you pull out the next ball. What is the probability that  $y$  of these balls are white? I.e., we are asking here for the counterpart of formula (4.2.1) if sampling is done with replacement.*

ANSWER.

$$(4.2.2) \quad \left(\frac{w}{n}\right)^y \left(\frac{n-w}{n}\right)^{m-y} \binom{m}{y}$$

□

Without proof we will state here that the expected value of  $\mathbf{y}$ , the number of white balls in the sample, is  $E[\mathbf{y}] = m\frac{w}{n}$ , which is the same as if one would select the balls with replacement.

Also without proof, the variance of  $\mathbf{y}$  is

$$(4.2.3) \quad \text{var}[\mathbf{y}] = m \frac{w}{n} \frac{(n-w)}{n} \frac{(n-m)}{(n-1)}.$$

This is smaller than the variance if one would choose with replacement, which is represented by the above formula without the last term  $\frac{n-m}{n-1}$ . This last term is called the finite population correction. More about all this is in [Lar82, p. 176–183].

### 4.3. The Poisson Distribution

The Poisson distribution counts the number of events in a given time interval. This number has the Poisson distribution if each event is the cumulative result of a large number of independent possibilities, each of which has only a small chance of occurring (law of rare events). The expected number of occurrences is proportional

to time with a proportionality factor  $\lambda$ , and in a short time span only zero or one event can occur, i.e., for infinitesimal time intervals it becomes a Bernoulli trial.

Approximate it by dividing the time from 0 to  $t$  into  $n$  intervals of length  $\frac{t}{n}$ ; then the occurrences are approximately  $n$  independent Bernoulli trials with probability of success  $\frac{\lambda t}{n}$ . (This is an approximation since some of these intervals may have more than one occurrence; but if the intervals become very short the probability of having two occurrences in the same interval becomes negligible.)

In this discrete approximation, the probability to have  $k$  successes in time  $t$  is

$$(4.3.1) \quad \Pr[\mathbf{x}=k] = \binom{n}{k} \left(\frac{\lambda t}{n}\right)^k \left(1 - \frac{\lambda t}{n}\right)^{(n-k)}$$

$$(4.3.2) \quad = \frac{1}{k!} \frac{n(n-1)\cdots(n-k+1)}{n^k} (\lambda t)^k \left(1 - \frac{\lambda t}{n}\right)^n \left(1 - \frac{\lambda t}{n}\right)^{-k}$$

$$(4.3.3) \quad \rightarrow \frac{(\lambda t)^k}{k!} e^{-\lambda t} \quad \text{for } n \rightarrow \infty \text{ while } k \text{ remains constant}$$

(4.3.3) is the limit because the second and the last term in (4.3.2)  $\rightarrow 1$ . The sum of all probabilities is 1 since  $\sum_{k=0}^{\infty} \frac{(\lambda t)^k}{k!} = e^{\lambda t}$ . The expected value is (note that we can have the sum start at  $k = 1$ ):

$$(4.3.4) \quad E[\mathbf{x}] = e^{-\lambda t} \sum_{k=1}^{\infty} k \frac{(\lambda t)^k}{k!} = \lambda t e^{-\lambda t} \sum_{k=1}^{\infty} \frac{(\lambda t)^{k-1}}{(k-1)!} = \lambda t.$$

This is the same as the expected value of the discrete approximations.

**PROBLEM 81.**  $\mathbf{x}$  follows a Poisson distribution, i.e.,

$$(4.3.5) \quad \Pr[\mathbf{x}=k] = \frac{(\lambda t)^k}{k!} e^{-\lambda t} \quad \text{for } k = 0, 1, \dots$$

• a. 2 points Show that  $E[\mathbf{x}] = \lambda t$ .

ANSWER. See (4.3.4).

• b. 4 points Compute  $E[\mathbf{x}(\mathbf{x} - 1)]$  and show that  $\text{var}[\mathbf{x}] = \lambda t$ .

ANSWER. For  $E[\mathbf{x}(\mathbf{x} - 1)]$  we can have the sum start at  $k = 2$ :

$$(4.3.6) \quad E[\mathbf{x}(\mathbf{x} - 1)] = e^{-\lambda t} \sum_{k=2}^{\infty} k(k-1) \frac{(\lambda t)^k}{k!} = (\lambda t)^2 e^{-\lambda t} \sum_{k=2}^{\infty} \frac{(\lambda t)^{k-2}}{(k-2)!} = (\lambda t)^2.$$

From this follows

$$(4.3.7) \quad \text{var}[\mathbf{x}] = E[\mathbf{x}^2] - (E[\mathbf{x}])^2 = E[\mathbf{x}(\mathbf{x} - 1)] + E[\mathbf{x}] - (E[\mathbf{x}])^2 = (\lambda t)^2 + \lambda t - (\lambda t)^2 = \lambda t.$$

The Poisson distribution can be used as an approximation to the Binomial distribution when  $n$  large,  $p$  small, and  $np$  moderate.



PROBLEM 82. Which value of  $\lambda$  would one need to approximate a given Binomial with  $n$  and  $p$ ?

ANSWER. That which gives the right expected value, i.e.,  $\lambda = np$ .  $\square$

PROBLEM 83. Two researchers counted cars coming down a road, which obey a Poisson distribution with unknown parameter  $\lambda$ . In other words, in an interval of length  $t$  one will have  $k$  cars with probability

$$(4.3.8) \quad \frac{(\lambda t)^k}{k!} e^{-\lambda t}.$$

Their assignment was to count how many cars came in the first half hour, and how many cars came in the second half hour. However they forgot to keep track of the time when the first half hour was over, and therefore wound up only with one count, namely, they knew that 213 cars had come down the road during this hour. They were afraid they would get fired if they came back with one number only, so they applied the following remedy: they threw a coin 213 times and counted the number of heads. This number, they pretended, was the number of cars in the first half hour.

• a. 6 points Did the probability distribution of the number gained in this way differ from the distribution of actually counting the number of cars in the first half hour?

ANSWER. First a few definitions:  $x$  is the total number of occurrences in the interval  $[0, 1]$ .  $y$  is the number of occurrences in the interval  $[0, t]$  (for a fixed  $t$ ; in the problem it was  $t = \frac{1}{2}$ , but we will do it for general  $t$ , which will make the notation clearer and more compact. Then we want to compute  $\Pr[y=m|x=n]$ . By definition of conditional probability:

$$(4.3.9) \quad \Pr[y=m|x=n] = \frac{\Pr[y=m \text{ and } x=n]}{\Pr[x=n]}.$$

How can we compute the probability of the intersection  $\Pr[y=m \text{ and } x=n]$ ? Use a trick: express this intersection as the intersection of independent events. For this define  $z$  as the number of events in the interval  $(t, 1]$ . Then  $\{y=m \text{ and } x=n\} = \{y=m \text{ and } z=n-m\}$ ; therefore  $\Pr[y=m \text{ and } x=n] = \Pr[y=m] \Pr[z=n-m]$ ; use this to get

$$(4.3.10) \quad \Pr[y=m|x=n] = \frac{\Pr[y=m] \Pr[z=n-m]}{\Pr[x=n]} = \frac{\frac{\lambda^m t^m}{m!} e^{-\lambda t} \frac{\lambda^{n-m} (1-t)^{n-m}}{(n-m)!} e^{-\lambda(1-t)}}{\frac{\lambda^n}{n!} e^{-\lambda}} = \binom{n}{m} t^m (1-t)^{n-m},$$

Here we use the fact that  $\Pr[x=k] = \frac{t^k}{k!} e^{-t}$ ,  $\Pr[y=k] = \frac{(\lambda t)^k}{k!} e^{-\lambda t}$ ,  $\Pr[z=k] = \frac{(1-\lambda)^k t^k}{k!} e^{-(1-\lambda)t}$ . One sees that a.  $\Pr[y=m|x=n]$  does not depend on  $\lambda$ , and b. it is exactly the probability of having  $m$  successes and  $n-m$  failures in a Bernoulli trial with success probability  $t$ . Therefore the procedure with the coins gave the two researchers a result which had the same probability distribution as if they had counted the number of cars in each half hour separately.  $\square$

• b. 2 points Explain what it means that the probability distribution of the number for the first half hour gained by throwing the coins does not differ from the one gained by actually counting the cars. Which condition is absolutely necessary for this hold?

ANSWER. The supervisor would never be able to find out through statistical analysis of data they delivered, even if they did it repeatedly. All estimation results based on the faked statistics would be as accurate regarding  $\lambda$  as the true statistics. All this is only true under the assumption that the cars really obey a Poisson distribution and that the coin is fair.

The fact that the Poisson as well as the binomial distributions are memoryless has nothing to do with them having a sufficient statistic.

PROBLEM 84. 8 points  $x$  is the number of customers arriving at a service counter in one hour.  $x$  follows a Poisson distribution with parameter  $\lambda = 2$ , i.e.,

$$(4.3.11) \quad \Pr[x=j] = \frac{2^j}{j!} e^{-2}.$$

• a. Compute the probability that only one customer shows up at the service counter during the hour, the probability that two show up, and the probability that one shows up.

• b. Despite the small number of customers, two employees are assigned to the service counter. They are hiding in the back, and whenever a customer steps up to the counter and rings the bell, they toss a coin. If the coin shows head, Herbert serves the customer, and if it shows tails, Karl does. Compute the probability that Herbert has to serve exactly one customer during the hour. Hint:

$$(4.3.12) \quad e = 1 + 1 + \frac{1}{2!} + \frac{1}{3!} + \frac{1}{4!} + \dots$$

• c. For any integer  $k \geq 0$ , compute the probability that Herbert has to serve exactly  $k$  customers during the hour.

PROBLEM 85. 3 points Compute the moment generating function of a Poisson variable observed over a unit time interval, i.e.,  $x$  satisfies  $\Pr[x=k] = \frac{\lambda^k}{k!} e^{-\lambda}$ . Do you want  $E[e^{tx}]$  for all  $t$ .

$$\text{ANSWER. } E[e^{tx}] = \sum_{k=0}^{\infty} e^{tk} \frac{\lambda^k}{k!} e^{-\lambda} = \sum_{k=0}^{\infty} \frac{(\lambda e^t)^k}{k!} e^{-\lambda} = e^{\lambda e^t} e^{-\lambda} = e^{\lambda(e^t-1)}.$$

#### 4.4. The Exponential Distribution

Now we will discuss random variables which are related to the Poisson distribution. At time  $t = 0$  you start observing a Poisson process, and the random variable  $t$  denotes the time you have to wait until the first occurrence.  $t$  can have any nonnegative real number as value. One can derive its cumulative distribution

as follows.  $t > t$  if and only if there are no occurrences in the interval  $[0, t]$ . Therefore  $\Pr[t > t] = \frac{(\lambda t)^0}{0!} e^{-\lambda t} = e^{-\lambda t}$ , and hence the cumulative distribution function  $F_t(t) = \Pr[t \leq t] = 1 - e^{-\lambda t}$  when  $t \geq 0$ , and  $F_t(t) = 0$  for  $t < 0$ . The density function is therefore  $f_t(t) = \lambda e^{-\lambda t}$  for  $t \geq 0$ , and 0 otherwise. This is called the exponential density function (its discrete analog is the geometric random variable). It can also be called a Gamma variable with parameters  $r = 1$  and  $\lambda$ .

**PROBLEM 86.** *2 points An exponential random variable  $t$  with parameter  $\lambda > 0$  has the density  $f_t(t) = \lambda e^{-\lambda t}$  for  $t \geq 0$ , and 0 for  $t < 0$ . Use this density to compute the expected value of  $t$ .*

**ANSWER.**  $E[t] = \int_0^\infty \lambda t e^{-\lambda t} dt = \int_0^\infty uv' dt = uv \Big|_0^\infty - \int_0^\infty u'v dt$ , where  $\frac{u=t}{u'=1} \frac{v'=\lambda e^{-\lambda t}}{v=-e^{-\lambda t}}$ . One can also use the more abbreviated notation  $= \int_0^\infty u dv = uv \Big|_0^\infty - \int_0^\infty v du$ , where  $\frac{u=t}{du'=dt} \frac{dv'=\lambda e^{-\lambda t}}{v=-e^{-\lambda t}} dt$ . Either way one obtains  $E[t] = -te^{-\lambda t} \Big|_0^\infty + \int_0^\infty e^{-\lambda t} dt = 0 - \frac{1}{\lambda} e^{-\lambda t} \Big|_0^\infty = \frac{1}{\lambda}$ .  $\square$

**PROBLEM 87.** *4 points An exponential random variable  $t$  with parameter  $\lambda > 0$  has the density  $f_t(t) = \lambda e^{-\lambda t}$  for  $t \geq 0$ , and 0 for  $t < 0$ . Use this density to compute the expected value of  $t^2$ .*

**ANSWER.** One can use that  $\Gamma(r) = \int_0^\infty \lambda^r t^{r-1} e^{-\lambda t} dt$  for  $r = 3$  to get:  $E[t^2] = (1/\lambda^2)\Gamma(3) = 2/\lambda^2$ . Or all from scratch:  $E[t^2] = \int_0^\infty \lambda t^2 e^{-\lambda t} dt = \int_0^\infty uv' dt = uv \Big|_0^\infty - \int_0^\infty u'v dt$ , where  $\frac{u=t^2}{u'=2t} \frac{v'=\lambda e^{-\lambda t}}{v=-e^{-\lambda t}}$ . Therefore  $E[t^2] = -t^2 e^{-\lambda t} \Big|_0^\infty + \int_0^\infty 2te^{-\lambda t} dt$ . The first term vanishes, for the second do it again:  $\int_0^\infty 2te^{-\lambda t} dt = \int_0^\infty uv' dt = uv \Big|_0^\infty - \int_0^\infty u'v dt$ , where  $\frac{u=t}{u'=1} \frac{v'=e^{-\lambda t}}{v=-(1/\lambda)e^{-\lambda t}}$ . Therefore the second term becomes  $2(t/\lambda)e^{-\lambda t} \Big|_0^\infty + 2 \int_0^\infty (1/\lambda)e^{-\lambda t} dt = 2/\lambda^2$ .  $\square$

**PROBLEM 88.** *2 points Does the exponential random variable with parameter  $\lambda > 0$ , whose cumulative distribution function is  $F_t(t) = 1 - e^{-\lambda t}$  for  $t \geq 0$ , and 0 otherwise, have a memory-less property? Compare Problem 76. Formulate this memory-less property and then verify whether it holds or not.*

**ANSWER.** Here is the formulation: for  $s < t$  follows  $\Pr[t > t | t > s] = \Pr[t > t - s]$ . This does indeed hold. Proof: lhs  $= \frac{\Pr[t > t \text{ and } t > s]}{\Pr[t > s]} = \frac{\Pr[t > t]}{\Pr[t > s]} = \frac{e^{-\lambda t}}{e^{-\lambda s}} = e^{-\lambda(t-s)}$ .  $\square$

**PROBLEM 89.** *The random variable  $t$  denotes the duration of an unemployment spell. It has the exponential distribution, which can be defined by:  $\Pr[t > t] = e^{-\lambda t}$  for  $t \geq 0$  ( $t$  cannot assume negative values).*

• a. *1 point Use this formula to compute the cumulative distribution function  $F_t(t)$  and the density function  $f_t(t)$*

**ANSWER.**  $F_t(t) = \Pr[t \leq t] = 1 - \Pr[t > t] = 1 - e^{-\lambda t}$  for  $t \geq 0$ , zero otherwise. Taking the derivative gives  $f_t(t) = \lambda e^{-\lambda t}$  for  $t \geq 0$ , zero otherwise.  $\square$

• b. *2 points What is the probability that an unemployment spell ends after time  $t + h$ , given that it has not yet ended at time  $t$ ? Show that this is the same as unconditional probability that an unemployment spell ends after time  $h$  (memory-less property).*

**ANSWER.**

$$(4.4.1) \quad \Pr[t > t + h | t > t] = \frac{\Pr[t > t + h]}{\Pr[t > t]} = \frac{e^{-\lambda(t+h)}}{e^{-\lambda t}} = e^{-\lambda h}$$

• c. *3 points Let  $h$  be a small number. What is the probability that an unemployment spell ends at or before  $t + h$ , given that it has not yet ended at time  $t$ ? How for small  $h$ , one can write approximately*

$$(4.4.2) \quad \Pr[t < t \leq t + h] = h f_t(t).$$

**ANSWER.**

$$(4.4.3) \quad \Pr[t \leq t + h | t > t] = \frac{\Pr[t \leq t + h \text{ and } t > t]}{\Pr[t > t]} = \frac{h f_t(t)}{1 - F_t(t)} = \frac{h \lambda e^{-\lambda t}}{e^{-\lambda t}} = h \lambda.$$

## 4.5. The Gamma Distribution

The time until the *second* occurrence of a Poisson event is a random variable which we will call  $t^{(2)}$ . Its cumulative distribution function is  $F_{t^{(2)}}(t) = \Pr[t^{(2)} \leq t] = 1 - \Pr[t^{(2)} > t]$ . But  $t^{(2)} > t$  means: there are either zero or one occurrences in the time between 0 and  $t$ ; therefore  $\Pr[t^{(2)} > t] = \Pr[x=0] + \Pr[x=1] = e^{-\lambda t} + \lambda t e^{-\lambda t}$ . Putting all together gives  $F_{t^{(2)}}(t) = 1 - e^{-\lambda t} - \lambda t e^{-\lambda t}$ . In order to differentiate the cumulative distribution function we need the product rule of differentiation:  $(uv)' = u'v + uv'$ . This gives

$$(4.5.1) \quad f_{t^{(2)}}(t) = \lambda e^{-\lambda t} - \lambda e^{-\lambda t} + \lambda^2 t e^{-\lambda t} = \lambda^2 t e^{-\lambda t}.$$

**PROBLEM 90.** *3 points Compute the density function of  $t^{(3)}$ , the time of the third occurrence of a Poisson variable.*

**ANSWER.**

$$(4.5.2) \quad \Pr[t^{(3)} > t] = \Pr[x=0] + \Pr[x=1] + \Pr[x=2]$$

$$(4.5.3) \quad F_{t^{(3)}}(t) = \Pr[t^{(3)} \leq t] = 1 - (1 + \lambda t + \frac{\lambda^2}{2} t^2) e^{-\lambda t}$$

$$(4.5.4) \quad f_{t^{(3)}}(t) = \frac{\partial}{\partial t} F_{t^{(3)}}(t) = -\left(-\lambda(1 + \lambda t + \frac{\lambda^2}{2} t^2) + (\lambda + \lambda^2 t)\right) e^{-\lambda t} = \frac{\lambda^3}{2} t^2 e^{-\lambda t}.$$

□

If one asks for the  $r$ th occurrence, again all but the last term cancel in the differentiation, and one gets

$$(4.5.5) \quad f_{t^{(r)}}(t) = \frac{\lambda^r}{(r-1)!} t^{r-1} e^{-\lambda t}.$$

This density is called the Gamma density with parameters  $\lambda$  and  $r$ .

The following definite integral, which is defined for all  $r > 0$  and all  $\lambda > 0$  is called the Gamma function:

$$(4.5.6) \quad \Gamma(r) = \int_0^\infty \lambda^r t^{r-1} e^{-\lambda t} dt.$$

Although this integral cannot be expressed in a closed form, it is an important function in mathematics. It is a well behaved function interpolating the factorials in the sense that  $\Gamma(r) = (r-1)!$ .

PROBLEM 91. Show that  $\Gamma(r)$  as defined in (4.5.6) is independent of  $\lambda$ , i.e., instead of (4.5.6) one can also use the simpler equation

$$(4.5.7) \quad \Gamma(r) = \int_0^\infty t^{r-1} e^{-t} dt.$$

PROBLEM 92. 3 points Show by partial integration that the Gamma function satisfies  $\Gamma(r+1) = r\Gamma(r)$ .

ANSWER. Start with

$$(4.5.8) \quad \Gamma(r+1) = \int_0^\infty \lambda^{r+1} t^r e^{-\lambda t} dt$$

and integrate by parts:  $\int u'v dt = uv - \int uv' dt$  with  $u' = \lambda e^{-\lambda t}$  and  $v = \lambda^r t^r$ , therefore  $u = -e^{-\lambda t}$  and  $v' = r\lambda^r t^{r-1}$ :

$$(4.5.9) \quad \Gamma(r+1) = -\lambda^r t^r e^{-\lambda t} \Big|_0^\infty + \int_0^\infty r\lambda^r t^{r-1} e^{-\lambda t} dt = 0 + r\Gamma(r).$$

□

PROBLEM 93. Show that  $\Gamma(r) = (r-1)!$  for all natural numbers  $r = 1, 2, \dots$

ANSWER. Proof by induction. First verify that it holds for  $r = 1$ , i.e., that  $\Gamma(1) = 1$ :

$$(4.5.10) \quad \Gamma(1) = \int_0^\infty \lambda e^{-\lambda t} dt = -e^{-\lambda t} \Big|_0^\infty = 1$$

and then, assuming that  $\Gamma(r) = (r-1)!$  Problem 92 says that  $\Gamma(r+1) = r\Gamma(r) = r(r-1)! = r!$ . □

Without proof:  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ . This will be shown in Problem 141.

Therefore the following defines a density function, called the Gamma density with parameter  $r$  and  $\lambda$ , for all  $r > 0$  and  $\lambda > 0$ :

$$(4.5.11) \quad f(x) = \frac{\lambda^r}{\Gamma(r)} x^{r-1} e^{-\lambda x} \quad \text{for } x \geq 0, \quad 0 \text{ otherwise.}$$

The only application we have for it right now is: this is the distribution of the time one has to wait until the  $r$ th occurrence of a Poisson distribution with intensity  $\lambda$ . Later we will have other applications in which  $r$  is not an integer.

PROBLEM 94. 4 points Compute the moment generating function of the Gamma distribution.

ANSWER.

$$(4.5.12) \quad m_x(t) = \mathbb{E}[e^{tx}] = \int_0^\infty e^{tx} \frac{\lambda^r}{\Gamma(r)} x^{r-1} e^{-\lambda x} dx$$

$$(4.5.13) \quad = \frac{\lambda^r}{(\lambda-t)^r} \int_0^\infty \frac{(\lambda-t)^r x^{r-1}}{\Gamma(r)} e^{-(\lambda-t)x} dx$$

$$(4.5.14) \quad = \left( \frac{\lambda}{\lambda-t} \right)^r$$

since the integrand in (4.5.12) is the density function of a Gamma distribution with parameter  $r$  and  $\lambda-t$ .

PROBLEM 95. 2 points The density and moment generating functions of a Gamma variable  $x$  with parameters  $r > 0$  and  $\lambda > 0$  are

$$(4.5.15) \quad f_x(x) = \frac{\lambda^r}{\Gamma(r)} x^{r-1} e^{-\lambda x} \quad \text{for } x \geq 0, \quad 0 \text{ otherwise.}$$

$$(4.5.16) \quad m_x(t) = \left( \frac{\lambda}{\lambda-t} \right)^r.$$

Show the following: If  $x$  has a Gamma distribution with parameters  $r$  and  $\lambda$ , then  $v = x/\lambda$  has a Gamma distribution with parameters  $r$  and 1. You can prove this either using the transformation theorem for densities, or the moment-generating function.

ANSWER. Solution using density function: The random variable whose density we know is  $x$  has density  $f_x(x) = \frac{\lambda^r}{\Gamma(r)} x^{r-1} e^{-\lambda x}$ . If  $x = \lambda v$ , then  $\frac{dx}{dv} = \lambda$ , and the absolute value is also  $\lambda$ . Therefore the density of  $v$  is  $f_v(v) = \frac{\lambda^r}{\Gamma(r)} (\lambda v)^{r-1} e^{-\lambda(\lambda v)}$ . Solution using the mgf:

$$(4.5.17) \quad m_x(t) = \mathbb{E}[e^{tx}] = \left( \frac{\lambda}{\lambda-t} \right)^r$$

$$(4.5.18) \quad m_v(t) \mathbb{E}[e^{tv}] = \mathbb{E}[e^{(t/\lambda)x}] = \left( \frac{1}{1-(t/\lambda)} \right)^r = \left( \frac{\lambda}{\lambda-t} \right)^r$$

but this last expression can be recognized to be the mgf of a Gamma with  $r$  and  $\lambda$ .

**PROBLEM 96.** 2 points It  $x$  has a Gamma distribution with parameters  $r$  and  $\lambda$ , and  $y$  one with parameters  $p$  and  $\lambda$ , and both are independent, show that  $x + y$  has a Gamma distribution with parameters  $r + p$  and  $\lambda$  (reproductive property of the Gamma distribution.) You may use equation (4.5.14) without proof

ANSWER.

$$(4.5.19) \quad \left(\frac{\lambda}{\lambda-t}\right)^r \left(\frac{\lambda}{\lambda-t}\right)^p = \left(\frac{\lambda}{\lambda-t}\right)^{r+p}.$$

□

**PROBLEM 97.** Show that a Gamma variable  $x$  with parameters  $r$  and  $\lambda$  has expected value  $E[x] = r/\lambda$  and variance  $\text{var}[x] = r/\lambda^2$ .

ANSWER. Proof with moment generating function:

$$(4.5.20) \quad \frac{d}{dt} \left(\frac{\lambda}{\lambda-t}\right)^r = \frac{r}{\lambda} \left(\frac{\lambda}{\lambda-t}\right)^{r+1},$$

therefore  $E[x] = \frac{r}{\lambda}$ , and by differentiating twice (apply the same formula again),  $E[x^2] = \frac{r(r+1)}{\lambda^2}$ , therefore  $\text{var}[x] = \frac{r}{\lambda^2}$ .

Proof using density function: For the expected value one gets  $E[t] = \int_0^\infty t \cdot \frac{\lambda^r}{\Gamma(r)} t^{r-1} e^{-\lambda t} dt = \frac{r}{\lambda} \frac{1}{\Gamma(r+1)} \int_0^\infty t^r \lambda^{r+1} e^{-\lambda t} dt = \frac{r}{\lambda} \frac{\Gamma(r+1)}{\Gamma(r+1)} = \frac{r}{\lambda}$ . Using the same tricks  $E[t^2] = \int_0^\infty t^2 \cdot \frac{\lambda^r}{\Gamma(r)} t^{r-1} e^{-\lambda t} dt = \frac{r(r+1)}{\lambda^2} \int_0^\infty \frac{\lambda^{r+2}}{\Gamma(r+2)} t^{r+1} e^{-\lambda t} dt = \frac{r(r+1)}{\lambda^2}$ .

Therefore  $\text{var}[t] = E[t^2] - (E[t])^2 = r/\lambda^2$ .

□

## 4.6. The Uniform Distribution

**PROBLEM 98.** Let  $x$  be uniformly distributed in the interval  $[a, b]$ , i.e., the density function of  $x$  is a constant for  $a \leq x \leq b$ , and zero otherwise.

- a. 1 point What is the value of this constant?

ANSWER. It is  $\frac{1}{b-a}$

□

- b. 2 points Compute  $E[x]$

ANSWER.  $E[x] = \int_a^b \frac{x}{b-a} dx = \frac{1}{b-a} \frac{b^2-a^2}{2} = \frac{a+b}{2}$  since  $b^2 - a^2 = (b+a)(b-a)$ .

□

- c. 2 points Show that  $E[x^2] = \frac{a^2+ab+b^2}{3}$ .

ANSWER.  $E[x^2] = \int_a^b \frac{x^2}{b-a} dx = \frac{1}{b-a} \frac{b^3-a^3}{3}$ . Now use the identity  $b^3 - a^3 = (b-a)(b^2 + ab + a^2)$  (check it by multiplying out).

□

- d. 2 points Show that  $\text{var}[x] = \frac{(b-a)^2}{12}$ .

ANSWER.  $\text{var}[x] = E[x^2] - (E[x])^2 = \frac{a^2+ab+b^2}{3} - \left(\frac{a+b}{2}\right)^2 = \frac{4a^2+4ab+4b^2}{12} - \frac{3a^2+6ab+3b^2}{12} = \frac{(b-a)^2}{12}$ .

□

## 4.7. The Beta Distribution

Assume you have two independent variables, both distributed uniformly over the interval  $[0, 1]$ , and you want to know the distribution of their maximum. Or their minimum. Or you have three and you want the distribution of the one in the middle. Then the densities have their maximum to the right, or to the left, or in the middle. The distribution of the  $r$ th highest out of  $n$  independent uniform variables is an example of the Beta density function. Can also be done and is probably theoretically meaningful for arbitrary real  $r$  and  $n$ .

**PROBLEM 99.**  $x$  and  $y$  are two independent random variables distributed uniformly over the interval  $[0, 1]$ . Let  $u$  be their minimum  $u = \min(x, y)$  (i.e.,  $u$  takes the value of  $x$  when  $x$  is smaller, and the value of  $y$  when  $y$  is smaller), and  $v = \max(x, y)$ .

- a. 2 points Given two numbers  $q$  and  $r$  between 0 and 1. Draw the events  $u > q$  and  $v \leq r$  into the unit square and compute their probabilities.

- b. 2 points Compute the density functions  $f_u(u)$  and  $f_v(v)$ .

- c. 2 points Compute the expected values of  $u$  and  $v$ .

ANSWER. For  $u$ :  $\Pr[u \leq q] = 1 - \Pr[u > q] = 1 - (1-q)^2 = 2q - q^2$ .  $f_u(v) = 2v$  Therefore  $f_u(u) = 2 - 2u$

$$(4.7.1) \quad E[u] = \int_0^1 (2-2u)u du = \left(u^2 - \frac{2u^3}{3}\right) \Big|_0^1 = \frac{1}{3}.$$

For  $v$  it is:  $\Pr[v \leq r] = r^2$ ; this is at the same time the cumulative distribution function. Therefore the density function is  $f_v(v) = 2v$  for  $0 \leq v \leq 1$  and 0 elsewhere.

$$(4.7.2) \quad E[v] = \int_0^1 v2v dv = \frac{2v^3}{3} \Big|_0^1 = \frac{2}{3}.$$

## 4.8. The Normal Distribution

By definition,  $y$  is normally distributed with mean  $\mu$  and variance  $\sigma^2$ , in symbols  $y \sim N(\mu, \sigma^2)$ , if it has the density function

$$(4.8.1) \quad f_y(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}.$$

It will be shown a little later that this is indeed a density function. This distribution has the highest entropy among all distributions with a given mean and variance. [Kap89, p. 47].

If  $y \sim N(\mu, \sigma^2)$ , then  $z = (y - \mu)/\sigma \sim N(0, 1)$ , which is called the standard Normal distribution.

PROBLEM 100. 2 points Compare [Gre97, p. 68]: Assume  $x \sim N(3, 4)$  (mean is 3 and variance 4). Determine with the help of a table of the Standard Normal Distribution function  $\Pr[2 < x \leq 5]$ .

ANSWER.  $\Pr[2 < x \leq 5] = \Pr[2 - 3 < x - 3 \leq 5 - 3] = \Pr[\frac{2-3}{2} < \frac{x-3}{2} \leq \frac{5-3}{2}] = \Pr[-\frac{1}{2} < \frac{x-3}{2} \leq 1] = \Phi(1) - \Phi(-\frac{1}{2}) = \Phi(1) - (1 - \Phi(\frac{1}{2})) = \Phi(1) + \Phi(\frac{1}{2}) - 1 = 0.8413 + 0.6915 - 1 = 0.5328$ . Some tables (Greene) give the area between 0 and all positive values; in this case it is  $0.3413 + 0.1915$ .  $\square$

The moment generating function of a standard normal  $z \sim N(0, 1)$  is the following integral:

$$(4.8.2) \quad m_z(t) = E[e^{tz}] = \int_{-\infty}^{+\infty} e^{tz} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz.$$

To solve this integral, complete the square in the exponent:

$$(4.8.3) \quad tz - \frac{z^2}{2} = \frac{t^2}{2} - \frac{1}{2}(z - t)^2;$$

Note that the first summand,  $\frac{t^2}{2}$ , no longer depends on  $z$ ; therefore the factor  $e^{\frac{t^2}{2}}$  can be written in front of the integral:

$$(4.8.4) \quad m_z(t) = e^{\frac{t^2}{2}} \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z-t)^2} dz = e^{\frac{t^2}{2}},$$

because now the integrand is simply the density function of a  $N(t, 1)$ .

A general univariate normal  $x \sim N(\mu, \sigma^2)$  can be written as  $x = \mu + \sigma z$  with  $z \sim N(0, 1)$ , therefore

$$(4.8.5) \quad m_x(t) = E[e^{(\mu+\sigma z)t}] = e^{\mu t} E[e^{\sigma z t}] = e^{(\mu + \sigma^2 t^2 / 2)}.$$

PROBLEM 101. Given two independent normal variables  $x \sim N(\mu_x, \sigma_x^2)$  and  $y \sim N(\mu_y, \sigma_y^2)$ . Using the moment generating function, show that

$$(4.8.6) \quad \alpha x + \beta y \sim N(\alpha \mu_x + \beta \mu_y, \alpha^2 \sigma_x^2 + \beta^2 \sigma_y^2).$$

ANSWER. Because of independence, the moment generating function of  $\alpha x + \beta y$  is the product of the m.g.f. of  $\alpha x$  and the one of  $\beta y$ :

$$(4.8.7) \quad m_{\alpha x + \beta y}(t) = e^{\mu_x \alpha t + \sigma_x^2 \alpha^2 t^2 / 2} e^{\mu_y \beta t + \sigma_y^2 \beta^2 t^2 / 2} = e^{(\mu_x \alpha + \mu_y \beta)t + (\sigma_x^2 \alpha^2 + \sigma_y^2 \beta^2)t^2 / 2},$$

which is the moment generating function of a  $N(\alpha \mu_x + \beta \mu_y, \alpha^2 \sigma_x^2 + \beta^2 \sigma_y^2)$ .  $\square$

We will say more about the univariate normal later when we discuss the multivariate normal distribution.

Sometimes it is also necessary to use the truncated normal distributions. If  $z$  is standard normal, then

$$(4.8.8) \quad E[z|z > z] = \frac{f_z(z)}{1 - F_z(z)}, \quad \text{var}[z|z > z] = 1 - \mu(\mu - z), \quad \text{where } \mu = E[z|z > z].$$

This expected value is therefore the ordinate of the density function at point  $z$  divided by the tail area of the tail over which  $z$  is known to vary. (This rule is only valid for the normal density function, not in general!) These kinds of results can be found in [JK70, pp. 81–83] or in the original paper [Coh50]

PROBLEM 102. Every customer entering a car dealership in a certain location can be thought of as having a reservation price  $y$  in his or her mind: if the car will be offered at or below this reservation price, then he or she will buy the car, otherwise there will be no sale. (Assume for the sake of the argument all cars are equally good. Assume this reservation price is Normally distributed with mean \$6000 and standard deviation \$1000 (if you randomly pick a customer and ask his or her reservation price). If a sale is made, a person's consumer surplus is the difference between the reservation price and the price actually paid, otherwise it is zero. For this question you will need the table for the standard normal cumulative distribution function.

• a. 2 points A customer is offered a car at a price of \$5800. The probability that he or she will take the car is .

ANSWER. We need  $\Pr[y \geq 5800]$ . If  $y=5800$  then  $z = \frac{y-6000}{1000} = -0.2$ ;  $\Pr[z \geq -0.2] = 1 - \Pr[z \leq -0.2] = 1 - 0.4207 = 0.5793$ .

• b. 3 points Since it is the 63rd birthday of the owner of the dealership, all cars in the dealership are sold for the price of \$6300. You pick at random one of the people coming out of the dealership. The probability that this person bought a car and his or her consumer surplus was more than \$500 is .

ANSWER. This is the unconditional probability that the reservation price was higher than \$6300 + \$500 = \$6800. i.e.,  $\Pr[y \geq 6800]$ . Define  $z = (y - \$6000) / \$1000$ . It is a standard normal,  $y \leq \$6800 \iff z \leq .8$ , Therefore  $p = 1 - \Pr[z \leq .8] = .2119$ .

• c. 4 points Here is an alternative scenario: Since it is the 63rd birthday of the owner of the dealership, all cars in the dealership are sold for the “birthday special” price of \$6300. You pick at random one of the people who bought one of these “birthday specials” priced \$6300. The probability that this person's consumer surplus was more than \$500 is .

The important part of this question is: it depends on the outcome of the experiment whether or not someone is included in the sample selection bias.

ANSWER. Here we need the conditional probability:

$$(4.8.9) \quad p = \Pr[y > \$6800 | y > \$6300] = \frac{\Pr[y > \$6800]}{\Pr[y > \$6300]} = \frac{1 - \Pr[y \leq \$6800]}{1 - \Pr[y \leq \$6300]}.$$

Again use the standard normal  $z = (y - \$6000)/\$1000$ . As before,  $y \leq \$6800 \iff z \leq .8$ , and  $y \leq \$6300 \iff z \leq .3$ . Therefore

$$(4.8.10) \quad p = \frac{1 - \Pr[z \leq .8]}{1 - \Pr[z \leq .3]} = \frac{.2119}{.3821} = .5546.$$

It depends on the layout of the normal distribution table how this should be looked up. □

• d. 5 points *We are still picking out customers that have bought the birthday specials. Compute the median value  $m$  of such a customer's consumer surplus. It is defined by*

$$(4.8.11) \quad \Pr[y > \$6300 + m | y > \$6300] = \Pr[y \leq \$6300 + m | y > \$6300] = 1/2.$$

ANSWER. Obviously,  $m \geq \$0$ . Therefore

$$(4.8.12) \quad \Pr[y > \$6300 + m | y > \$6300] = \frac{\Pr[y > \$6300 + m]}{\Pr[y > \$6300]} = \frac{1}{2},$$

or  $\Pr[y > \$6300 + m] = (1/2) \Pr[y > \$6300] = (1/2) \cdot .3821 = .1910$ . I.e.,  $\Pr[\frac{y-6000}{1000} > \frac{6300-6000+m}{1000}] = \frac{.300}{1000} + \frac{m}{1000} = .1910$ . For this we find in the table  $\frac{.300}{1000} + \frac{m}{1000} = 0.875$ , therefore  $300 + m = 875$ , or  $m = \$575$ . □

• e. 3 points *Is the expected value of the consumer surplus of all customers that have bought a birthday special larger or smaller than the median? Fill in your answer*

here: . *Proof is not required, as long as the answer is correct.*

ANSWER. The mean is larger because it is more heavily influenced by outliers.

$$(4.8.13) \quad E[y - 6300 | y \geq 6300] = E[6000 + 1000z - 6300 | 6000 + 1000z \geq 6300]$$

$$(4.8.14) \quad = E[1000z - 300 | 1000z \geq 300]$$

$$(4.8.15) \quad = E[1000z | z \geq 0.3] - 300$$

$$(4.8.16) \quad = 1000 E[z | z \geq 0.3] - 300$$

$$(4.8.17) \quad = 1000 \frac{f(0.3)}{1 - \Psi(0.3)} - 300 = 698 > 575. \quad \square$$

### 4.9. The Chi-Square Distribution

A  $\chi^2$  with *one* degree of freedom is defined to be the distribution of the square  $q = z^2$  of a univariate standard normal variable.

Call the cumulative distribution function of a standard normal  $F_z(z)$ . Then the cumulative distribution function of the  $\chi^2$  variable  $q = z^2$  is, according to Problem 47,  $F_q(q) = 2F_z(\sqrt{q}) - 1$ . To get the density of  $q$  take the derivative of  $F_q(q)$  with

respect to  $q$ . For this we need the chain rule, first taking the derivative with respect to  $z = \sqrt{q}$  and multiply by  $\frac{dz}{dq}$ :

$$(4.9.1) \quad f_q(q) = \frac{d}{dq} (2F_z(\sqrt{q}) - 1) = \frac{d}{dq} (2F_z(z) - 1)$$

$$(4.9.2) \quad = 2 \frac{dF_z}{dz}(z) \frac{dz}{dq} = \frac{2}{\sqrt{2\pi}} e^{-z^2/2} \frac{1}{2\sqrt{q}}$$

$$(4.9.3) \quad = \frac{1}{\sqrt{2\pi q}} e^{-q/2}.$$

Now remember the Gamma function. Since  $\Gamma(1/2) = \sqrt{\pi}$  (Proof in Problem 14) one can rewrite (4.9.3) as

$$(4.9.4) \quad f_q(q) = \frac{(1/2)^{1/2} q^{-1/2} e^{-q/2}}{\Gamma(1/2)},$$

i.e., it is a Gamma density with parameters  $r = 1/2$ ,  $\lambda = 1/2$ .

A  $\chi^2$  with  $p$  degrees of freedom is defined as the sum of  $p$  independent univariate  $\chi^2$  variables. By the reproductive property of the Gamma distribution (Problem 9) this gives a Gamma variable with parameters  $r = p/2$  and  $\lambda = 1/2$ .

$$(4.9.5) \quad \text{If } q \sim \chi_p^2 \quad \text{then } E[q] = p \quad \text{and } \text{var}[q] = 2p$$

We will say that a random variable  $q$  is distributed as a  $\sigma^2 \chi_p^2$  iff  $q/\sigma^2$  is a  $\chi_p^2$ . This is the distribution of a sum of  $p$  independent  $N(0, \sigma^2)$  variables.

### 4.10. The Lognormal Distribution

This is a random variable whose log has a normal distribution. See [Gre97, 71]. Parametrized by the  $\mu$  and  $\sigma^2$  of its log. Density is

$$(4.10.1) \quad \frac{1}{x\sqrt{2\pi\sigma^2}} e^{-(\ln x - \mu/\sigma^2)/2}$$

[Cow77, pp. 82–87] has an excellent discussion of the properties of the lognormal for income distributions.

### 4.11. The Cauchy Distribution

PROBLEM 103. 6 points [JK70, pp. 155/6] *An example of a distribution with the same mean and variance is the Cauchy distribution, whose density looks much like a normal density, but has much thicker tails. The density and characteristic function are (I am not asking you to compute the characteristic function)*

$$(4.11.1) \quad f_x(x) = \frac{1}{\pi(1+x^2)} \quad E[e^{itx}] = \exp(-|t|).$$

Here  $i = \sqrt{-1}$ , but you should not be afraid of it, in most respects,  $i$  behaves like any real number. The characteristic function has properties very similar to the moment generating function, with the added advantage that it always exists. Using the characteristic functions show that if  $\mathbf{x}$  and  $\mathbf{y}$  are independent Cauchy distributions, then  $(\mathbf{x} + \mathbf{y})/2$  has the same distribution as  $\mathbf{x}$  or  $\mathbf{y}$ .

ANSWER.

$$(4.11.2) \quad \mathbb{E} \left[ \exp \left( it \frac{\mathbf{x} + \mathbf{y}}{2} \right) \right] = \mathbb{E} \left[ \exp \left( i \frac{t}{2} \mathbf{x} \right) \exp \left( i \frac{t}{2} \mathbf{y} \right) \right] = \exp \left( - \left| \frac{t}{2} \right| \right) \exp \left( - \left| \frac{t}{2} \right| \right) = \exp(-|t|).$$

□

It has taken a historical learning process to distinguish significant from insignificant events. The order in which the birds sit down on a tree is insignificant, but the constellation of stars on the night sky is highly significant for the seasons etc. The confusion between significant and insignificant events can explain how astrology arose: after it was discovered that the constellation of stars was significant, but without knowledge of the mechanism through which the constellation of stars was significant, people experimented to find evidence of causality between those aspects of the night sky that were changing, like the locations of the planets, and events on earth, like the births of babies. Romans thought the constellation of birds in the sky was significant.

Freud discovered that human error may be significant. Modern political consciousness still underestimates the extent to which the actions of states are significant: If a welfare recipient is faced with an intractable labyrinth of regulations and a multitude of agencies, then this is not the unintended result of bureaucracy gone wild, but it is deliberate: this bureaucratic nightmare deters people from using welfare, but it creates the illusion that welfare exists and it does give relief in some blatant cases.

Also “mistakes” like the bombing of the Chinese embassy are not mistakes but are significant.

In statistics the common consensus is that the averages are significant and the deviations from the averages are insignificant. By taking averages one distills the significant, systematic part of the data from the insignificant part. Usually this is justified by the “law of large numbers.” I.e., people think that this is something about reality which can be derived and proved mathematically. However this is an unrealistic position: how can math tell us which events are significant?

Here the Cauchy distribution is an interesting counterexample: it is a probability distribution for which it does not make sense to take averages. If one takes the average of  $n$  observations, then this average does not have less randomness than each individual observation, but it has exactly the same distribution as one single

observation. (The law of large numbers does not apply here because the Cauchy distribution does not have an expected value.)

In a world in which random outcomes are Cauchy-distributed, taking averages is not a good way to learn from one’s experiences. People who try to keep track of things by taking averages (or by running regressions, which is a natural extension of taking averages) would have the same status in that world as astrologers have in our world. Taking medians and other quantiles would be considered scientific, but taking averages would be considered superstition.

The lesson of this is: even a scientific procedure as innocuous as that of taking averages cannot be justified on purely epistemological grounds. Although it is widely assumed that the law of large numbers is such a justification, it is not. The law of large numbers does not always hold; it only holds if the random variable under consideration has an expected value.

The transcendental realist can therefore say: since it apparently does make sense to take averages in our world, we can deduce transcendently that many random variables which we are dealing with do have finite expected values.

This is perhaps the simplest case of a transcendental conclusion. But this simplest case also vindicates another one of Bhaskar’s assumptions: these transcendental conclusions cannot be arrived at in a non-transcendental way, by staying in the scientific sense itself. It is impossible to decide, using statistical means alone, whether our data come from a distribution which has finite expected values or not. The reason is that one always has only finite datasets, and the empirical distribution of a finite sample always has finite expected values, even if the sample comes from a population which does not have finite expected values.

## CHAPTER 5

## Chebyshev Inequality, Weak Law of Large Numbers, and Central Limit Theorem

### 5.1. Chebyshev Inequality

If the random variable  $\mathbf{y}$  has finite expected value  $\mu$  and standard deviation  $\sigma$ , and  $k$  is some positive number, then the *Chebyshev Inequality* says

$$(5.1.1) \quad \Pr[|\mathbf{y} - \mu| \geq k\sigma] \leq \frac{1}{k^2}.$$

In words, the probability that a given random variable  $\mathbf{y}$  differs from its expected value by more than  $k$  standard deviations is less than  $1/k^2$ . (Here “more than” and “less than” are short forms for “more than or equal to” and “less than or equal to.”) One does not need to know the full distribution of  $\mathbf{y}$  for that, only its expected value and standard deviation. We will give here a proof only if  $\mathbf{y}$  has a discrete distribution, but the inequality is valid in general. Going over to the standardized variable  $z = \frac{\mathbf{y} - \mu}{\sigma}$  we have to show  $\Pr[|z| \geq k] \leq \frac{1}{k^2}$ . Assuming  $z$  assumes the values  $z_1, z_2, \dots$  with probabilities  $p(z_1), p(z_2), \dots$ , then

$$(5.1.2) \quad \Pr[|z| \geq k] = \sum_{i: |z_i| \geq k} p(z_i).$$

Now multiply by  $k^2$ :

$$(5.1.3) \quad k^2 \Pr[|z| \geq k] = \sum_{i: |z_i| \geq k} k^2 p(z_i)$$

$$(5.1.4) \quad \leq \sum_{i: |z_i| \geq k} z_i^2 p(z_i)$$

$$(5.1.5) \quad \leq \sum_{\text{all } i} z_i^2 p(z_i) = \text{var}[z] = 1.$$

The Chebyshev inequality is *sharp* for all  $k \geq 1$ . Proof: the random variable which takes the value  $-k$  with probability  $\frac{1}{2k^2}$  and the value  $+k$  with probability

$\frac{1}{2k^2}$ , and 0 with probability  $1 - \frac{1}{k^2}$ , has expected value 0 and variance 1 and the  $\leq$ -sign in (5.1.1) becomes an equal sign.

PROBLEM 104. [HT83, p. 316] Let  $\mathbf{y}$  be the number of successes in  $n$  trials of a Bernoulli experiment with success probability  $p$ . Show that

$$(5.1.6) \quad \Pr\left(\left|\frac{\mathbf{y}}{n} - p\right| < \varepsilon\right) \geq 1 - \frac{1}{4n\varepsilon^2}.$$

*Hint: first compute what Chebyshev will tell you about the lefthand side, and then you will need still another inequality.*

ANSWER.  $E[\mathbf{y}/n] = p$  and  $\text{var}[\mathbf{y}/n] = pq/n$  (where  $q = 1 - p$ ). Chebyshev says therefore

$$(5.1.7) \quad \Pr\left(\left|\frac{\mathbf{y}}{n} - p\right| \geq k\sqrt{\frac{pq}{n}}\right) \leq \frac{1}{k^2}.$$

Setting  $\varepsilon = k\sqrt{pq/n}$ , therefore  $1/k^2 = pq/n\varepsilon^2$  one can rewrite (5.1.7) as

$$(5.1.8) \quad \Pr\left(\left|\frac{\mathbf{y}}{n} - p\right| \geq \varepsilon\right) \leq \frac{pq}{n\varepsilon^2}.$$

Now note that  $pq \leq 1/4$  whatever their values are.

PROBLEM 105. 2 points For a standard normal variable,  $\Pr[|z| \geq 1]$  is approximately 1/3, please look up the precise value in a table. What does the Chebyshev inequality say about this probability? Also,  $\Pr[|z| \geq 2]$  is approximately 5%, again look up the precise value. What does Chebyshev say?

ANSWER.  $\Pr[|z| \geq 1] = 0.3174$ , the Chebyshev inequality says that  $\Pr[|z| \geq 1] \leq 1$ . Also,  $\Pr[|z| \geq 2] = 0.0540$ , while Chebyshev says it is  $\leq 0.25$ .

### 5.2. The Probability Limit and the Law of Large Numbers

Let  $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \dots$  be a sequence of independent random variables all of which have the same expected value  $\mu$  and variance  $\sigma^2$ . Then  $\bar{\mathbf{y}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i$  has expected value  $\mu$  and variance  $\frac{\sigma^2}{n}$ . I.e., its probability mass is clustered much more closely around the value  $\mu$  than the individual  $\mathbf{y}_i$ . To make this statement more precise we need a concept of convergence of random variables. It is not possible to define it the “obvious” way that the sequence of random variables  $\mathbf{y}_n$  converges toward  $\mathbf{y}$  if every realization of them converges, since it is possible, although extremely unlikely, that e.g. all throws of a coin show heads ad infinitum, or follow another sequence for which the average number of heads does not converge towards 1/2. Therefore we will use the following definition:

The sequence of random variables  $\mathbf{y}_1, \mathbf{y}_2, \dots$  converges in probability to another random variable  $\mathbf{y}$  if and only if for every  $\delta > 0$

$$(5.2.1) \quad \lim_{n \rightarrow \infty} \Pr[|\mathbf{y}_n - \mathbf{y}| \geq \delta] = 0.$$



One can also say that the probability limit of  $\mathbf{y}_n$  is  $\mathbf{y}$ , in formulas

$$(5.2.2) \quad \text{plim}_{n \rightarrow \infty} \mathbf{y}_n = \mathbf{y}.$$

In many applications, the limiting variable  $\mathbf{y}$  is a degenerate random variable, i.e., it is a constant.

The *Weak Law of Large Numbers* says that, if the expected value exists, then the probability limit of the sample means of an ever increasing sample is the expected value, i.e.,  $\text{plim}_{n \rightarrow \infty} \bar{\mathbf{y}}_n = \mu$ .

**PROBLEM 106.** 5 points Assuming that not only the expected value but also the variance exists, derive the Weak Law of Large Numbers, which can be written as

$$(5.2.3) \quad \lim_{n \rightarrow \infty} \Pr[|\bar{\mathbf{y}}_n - \mathbf{E}[\mathbf{y}]| \geq \delta] = 0 \text{ for all } \delta > 0,$$

from the Chebyshev inequality

$$(5.2.4) \quad \Pr[|x - \mu| \geq k\sigma] \leq \frac{1}{k^2} \quad \text{where } \mu = \mathbf{E}[x] \text{ and } \sigma^2 = \text{var}[x]$$

**ANSWER.** From nonnegativity of probability and the Chebyshev inequality for  $x = \bar{\mathbf{y}}$  follows  $0 \leq \Pr[|\bar{\mathbf{y}} - \mu| \geq \frac{k\sigma}{\sqrt{n}}] \leq \frac{1}{k^2}$  for all  $k$ . Set  $k = \frac{\delta\sqrt{n}}{\sigma}$  to get  $0 \leq \Pr[|\bar{\mathbf{y}}_n - \mu| \geq \delta] \leq \frac{\sigma^2}{n\delta^2}$ . For any fixed  $\delta > 0$ , the upper bound converges towards zero as  $n \rightarrow \infty$ , and the lower bound is zero, therefore the probability itself also converges towards zero.  $\square$

**PROBLEM 107.** 4 points Let  $\mathbf{y}_1, \dots, \mathbf{y}_n$  be a sample from some unknown probability distribution, with sample mean  $\bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i$  and sample variance  $s^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})^2$ . Show that the data satisfy the following “sample equivalent” of the Chebyshev inequality: if  $k$  is any fixed positive number, and  $m$  is the number of observations  $\mathbf{y}_j$  which satisfy  $|\mathbf{y}_j - \bar{\mathbf{y}}| \geq ks$ , then  $m \leq n/k^2$ . In symbols,

$$(5.2.5) \quad \#\{\mathbf{y}_i : |\mathbf{y}_i - \bar{\mathbf{y}}| \geq ks\} \leq \frac{n}{k^2}.$$

*Hint: apply the usual Chebyshev inequality to the so-called empirical distribution of the sample. The empirical distribution is a discrete probability distribution defined by  $\Pr[\mathbf{y} = \mathbf{y}_i] = k/n$ , when the number  $\mathbf{y}_i$  appears  $k$  times in the sample. (If all  $\mathbf{y}_i$  are different, then all probabilities are  $1/n$ ). The empirical distribution corresponds to the experiment of randomly picking one observation out of the given sample.*

**ANSWER.** The only thing to note is: the sample mean is the expected value in that empirical distribution, the sample variance is the variance, and the relative number  $m/n$  is the probability.

$$(5.2.6) \quad \#\{\mathbf{y}_i : \mathbf{y}_i \in S\} = n \Pr[S] \quad \square$$

• a. 3 points What happens to this result when the distribution from which the  $\mathbf{y}_i$  are taken does not have an expected value or a variance?

**ANSWER.** The result still holds but  $\bar{\mathbf{y}}$  and  $s^2$  do not converge as the number of observations increases.

### 5.3. Central Limit Theorem

Assume all  $\mathbf{y}_i$  are independent and have the same distribution with mean  $\mu$  and variance  $\sigma^2$ , and also a moment generating function. Again, let  $\bar{\mathbf{y}}_n$  be the sample mean of the first  $n$  observations. The central limit theorem says that the probability distribution for

$$(5.3.1) \quad \frac{\bar{\mathbf{y}}_n - \mu}{\sigma/\sqrt{n}}$$

converges to a  $\mathbf{N}(0, 1)$ . This is a different concept of convergence than the probability limit, it is convergence in distribution.

**PROBLEM 108.** 1 point Construct a sequence of random variables  $\mathbf{y}_1, \mathbf{y}_2, \dots$  with the following property: their cumulative distribution functions converge to the cumulative distribution function of a standard normal, but the random variables themselves do not converge in probability. (This is easy!)

**ANSWER.** One example would be: all  $\mathbf{y}_i$  are independent standard normal variables.

Why do we have the funny expression  $\frac{\bar{\mathbf{y}}_n - \mu}{\sigma/\sqrt{n}}$ ? Because this is the standardized version of  $\bar{\mathbf{y}}_n$ . We know from the law of large numbers that the distribution of  $\bar{\mathbf{y}}_n$  becomes more and more concentrated around  $\mu$ . If we standardize the sample averages  $\bar{\mathbf{y}}_n$ , we compensate for this concentration. The central limit theorem tells us therefore what happens to the *shape* of the cumulative distribution function of  $\bar{\mathbf{y}}_n$ . If we disregard the fact that it becomes more and more concentrated (by multiplying it by a factor which is chosen such that the variance remains constant), then we see that its geometric shape comes closer and closer to a normal distribution.

Proof of the Central Limit Theorem: By Problem 109,

$$(5.3.2) \quad \frac{\bar{\mathbf{y}}_n - \mu}{\sigma/\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\mathbf{y}_i - \mu}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n z_i \quad \text{where } z_i = \frac{\mathbf{y}_i - \mu}{\sigma}.$$

Let  $m_3, m_4, \dots$ , be the third, fourth, etc., moments of  $z_i$ ; then the m.g.f. of  $z_i$  is

$$(5.3.3) \quad m_{z_i}(t) = 1 + \frac{t^2}{2!} + \frac{m_3 t^3}{3!} + \frac{m_4 t^4}{4!} + \dots$$

Therefore the m.g.f. of  $\frac{1}{\sqrt{n}} \sum_{i=1}^n z_i$  is (multiply and substitute  $t/\sqrt{n}$  for  $t$ ):

$$(5.3.4) \quad \left(1 + \frac{t^2}{2!n} + \frac{m_3 t^3}{3!\sqrt{n^3}} + \frac{m_4 t^4}{4!n^2} + \dots\right)^n = \left(1 + \frac{w_n}{n}\right)^n$$

where

$$(5.3.5) \quad w_n = \frac{t^2}{2!} + \frac{m_3 t^3}{3! \sqrt{n}} + \frac{m_4 t^4}{4! n} + \cdots .$$

Now use Euler's limit, this time in the form: if  $w_n \rightarrow w$  for  $n \rightarrow \infty$ , then  $\left(1 + \frac{w_n}{n}\right)^n \rightarrow e^w$ . Since our  $w_n \rightarrow \frac{t^2}{2}$ , the m.g.f. of the standardized  $\bar{y}_n$  converges toward  $e^{\frac{t^2}{2}}$ , which is that of a standard normal distribution.

The Central Limit theorem is an example of emergence: independently of the distributions of the individual summands, the distribution of the sum has a very specific shape, the Gaussian bell curve. The signals turn into white noise. Here emergence is the emergence of homogeneity and indeterminacy. In capitalism, much more specific outcomes emerge: whether one quits the job or not, whether one sells the stock or not, whether one gets a divorce or not, the outcome for society is to perpetuate the system. Not many activities don't have this outcome.

PROBLEM 109. Show in detail that  $\frac{\bar{y}_n - \mu}{\sigma/\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{y_i - \mu}{\sigma}$ .

ANSWER. Lhs =  $\frac{\sqrt{n}}{\sigma} \left( \left( \frac{1}{n} \sum_{i=1}^n y_i \right) - \mu \right) = \frac{\sqrt{n}}{\sigma} \left( \left( \frac{1}{n} \sum_{i=1}^n y_i \right) - \left( \frac{1}{n} \sum_{i=1}^n \mu \right) \right) = \frac{\sqrt{n}}{\sigma} \frac{1}{n} \left( \sum_{i=1}^n (y_i - \mu) \right) =$  rhs.  $\square$

PROBLEM 110. 3 points Explain verbally clearly what the law of large numbers means, what the Central Limit Theorem means, and what their difference is.

PROBLEM 111. (For this problem, a table is needed.) [Lar82, exercise 5.6.1, p. 301] If you roll a pair of dice 180 times, what is the approximate probability that the sum seven appears 25 or more times? Hint: use the Central Limit Theorem (but don't worry about the continuity correction, which is beyond the scope of this class).

ANSWER. Let  $x_i$  be the random variable that equals one if the  $i$ -th roll is a seven, and zero otherwise. Since 7 can be obtained in six ways (1+6, 2+5, 3+4, 4+3, 5+2, 6+1), the probability to get a 7 (which is at the same time the expected value of  $x_i$ ) is  $6/36=1/6$ . Since  $x_i^2 = x_i$ ,  $\text{var}[x_i] = E[x_i] - (E[x_i])^2 = \frac{1}{6} - \frac{1}{36} = \frac{5}{36}$ . Define  $x = \sum_{i=1}^{180} x_i$ . We need  $\Pr[x \geq 25]$ . Since  $x$  is the sum of many independent identically distributed random variables, the CLT says that  $x$  is asymptotically normal. Which normal? That which has the same expected value and variance as  $x$ .  $E[x] = 180 \cdot (1/6) = 30$  and  $\text{var}[x] = 180 \cdot (5/36) = 25$ . Therefore define  $y \sim N(30, 25)$ . The CLT says that  $\Pr[x \geq 25] \approx \Pr[y \geq 25]$ . Now  $y \geq 25 \iff y - 30 \geq -5 \iff y - 30 \leq +5 \iff (y - 30)/5 \leq 1$ . But  $z = (y - 30)/5$  is a standard Normal, therefore  $\Pr[(y - 30)/5 \leq 1] = F_z(1)$ , i.e., the cumulative distribution of the standard Normal evaluated at +1. One can look this up in a table, the probability asked for is .8413. Larson uses the continuity correction:  $x$  is discrete, and  $\Pr[x \geq 25] = \Pr[x > 24]$ . Therefore  $\Pr[y \geq 25]$  and  $\Pr[y > 24]$  are two alternative good approximations; but the best is  $\Pr[y \geq 24.5] = .8643$ . This is the continuity correction.  $\square$

CHAPTER 6

## Vector Random Variables

In this chapter we will look at *two* random variables  $\mathbf{x}$  and  $\mathbf{y}$  defined on the same sample space  $U$ , i.e.,

$$(6.0.6) \quad \mathbf{x}: U \ni \omega \mapsto x(\omega) \in \mathbb{R} \quad \text{and} \quad \mathbf{y}: U \ni \omega \mapsto y(\omega) \in \mathbb{R}.$$

As we said before,  $\mathbf{x}$  and  $\mathbf{y}$  are called independent if all events of the form  $x \leq x$  are independent of any event of the form  $y \leq y$ . But now let us assume they are *not* independent. In this case, we do not have all the information about them if we merely know the distribution of each.

The following example from [Lar82, example 5.1.7. on p. 233] illustrates the issues involved. This example involves two random variables that have only two possible outcomes each. Suppose you are told that a coin is to be flipped two times and that the probability of a head is .5 for each flip. This information is not enough to determine the probability of the second flip giving a head conditionally on the first flip giving a head.

For instance, the above two probabilities can be achieved by the following experimental setup: a person has one fair coin and flips it twice in a row. Then the two flips are independent.

But the probabilities of 1/2 for heads and 1/2 for tails can also be achieved as follows: The person has two coins in his or her pocket. One has two heads, and one has two tails. If at random one of these two coins is picked and flipped twice, then the second flip has the same outcome as the first flip.

What do we need to get the full picture? We must consider the two variables not separately but jointly, as a *totality*. In order to do this, we combine  $\mathbf{x}$  and  $\mathbf{y}$  into one entity, a vector  $\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in \mathbb{R}^2$ . Consequently we need to know the probability measure induced by the mapping  $U \ni \omega \mapsto \begin{bmatrix} x(\omega) \\ y(\omega) \end{bmatrix} \in \mathbb{R}^2$ .

It is not sufficient to look at random variables individually; one must look at them as a totality.

Therefore let us first get an overview over all possible probability measures on the plane  $\mathbb{R}^2$ . In strict analogy with the one-dimensional case, these probability measures

can be represented by the joint cumulative distribution function. It is defined as

$$(6.0.7) \quad F_{\mathbf{x},\mathbf{y}}(x, y) = \Pr\left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \leq \begin{bmatrix} x \\ y \end{bmatrix}\right] = \Pr[x \leq x \text{ and } \mathbf{y} \leq y].$$

For discrete random variables, for which the cumulative distribution function is a step function, the joint probability mass function provides the same information

$$(6.0.8) \quad p_{\mathbf{x},\mathbf{y}}(x, y) = \Pr\left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix}\right] = \Pr[x=x \text{ and } \mathbf{y}=y].$$

**PROBLEM 112.** Write down the joint probability mass functions for the two versions of the two coin flips discussed above.

**ANSWER.** Here are the probability mass functions for these two cases:

		Second Flip						Second Flip							
			<i>H</i>	<i>T</i>	sum		<i>H</i>	<i>T</i>	sum		<i>H</i>	<i>T</i>	sum		
(6.0.9)	First	<i>H</i>	.25	.25	.50	First	<i>H</i>	.50	.00	.50	First	<i>H</i>	.50	.00	.50
	Flip	<i>T</i>	.25	.25	.50		Flip	<i>T</i>	.00	.50	Flip	<i>T</i>	.00	.50	.50
		sum	.50	.50	1.00		sum	.50	.50	1.00		sum	.50	.50	1.00

The most important case is that with a differentiable cumulative distribution function. Then the joint density function  $f_{\mathbf{x},\mathbf{y}}(x, y)$  can be used to define the probability measure. One obtains it from the cumulative distribution function by taking derivatives:

$$(6.0.10) \quad f_{\mathbf{x},\mathbf{y}}(x, y) = \frac{\partial^2}{\partial x \partial y} F_{\mathbf{x},\mathbf{y}}(x, y).$$

Probabilities can be obtained back from the density function either by the integral condition, or by the infinitesimal condition. I.e., either one says for a subset  $B \subset \mathbb{R}^2$ :

$$(6.0.11) \quad \Pr\left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in B\right] = \int \int_B f(x, y) dx dy,$$

or one says, for a infinitesimal two-dimensional volume element  $dV_{x,y}$  located at  $(x, y)$  which has the two-dimensional volume (i.e., area)  $|dV|$ ,

$$(6.0.12) \quad \Pr\left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in dV_{x,y}\right] = f(x, y) |dV|.$$

The vertical bars here do not mean the absolute value but the volume of the argument inside.

### 6.1. Expected Value, Variances, Covariances

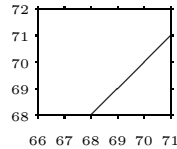
To get the expected value of a function of  $\mathbf{x}$  and  $\mathbf{y}$ , one simply has to put this function together with the density function into the integral, i.e., the formula is

$$(6.1.1) \quad E[g(\mathbf{x}, \mathbf{y})] = \int \int_{\mathbb{R}^2} g(x, y) f_{\mathbf{x}, \mathbf{y}}(x, y) dx dy.$$

PROBLEM 113. Assume there are two transportation choices available: bus and car. If you pick at random a neoclassical individual  $\omega$  and ask which utility this person derives from using bus or car, the answer will be two numbers that can be written as a vector  $\begin{bmatrix} u(\omega) \\ v(\omega) \end{bmatrix}$  ( $u$  for bus and  $v$  for car).

• a. 3 points Assuming  $\begin{bmatrix} u \\ v \end{bmatrix}$  has a uniform density in the rectangle with corners  $\begin{bmatrix} 66 \\ 68 \end{bmatrix}$ ,  $\begin{bmatrix} 66 \\ 72 \end{bmatrix}$ ,  $\begin{bmatrix} 71 \\ 68 \end{bmatrix}$ , and  $\begin{bmatrix} 71 \\ 72 \end{bmatrix}$ , compute the probability that the bus will be preferred.

ANSWER. The probability is 9/40.  $u$  and  $v$  have a joint density function that is uniform in the rectangle below and zero outside ( $u$ , the preference for buses, is on the horizontal, and  $v$ , the preference for cars, on the vertical axis). The probability is the fraction of this rectangle below the diagonal.



□

• b. 2 points How would you criticize an econometric study which argued along the above lines?

ANSWER. The preferences are not for a bus or a car, but for a whole transportation systems. And these preferences are not formed independently and individualistically, but they depend on which other infrastructures are in place, whether there is suburban sprawl or concentrated walkable cities, etc. This is again the error of detotalization (which favors the status quo).

□

Jointly distributed random variables should be written as random *vectors*. Instead of  $\begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix}$  we will also write  $\mathbf{x}$  (bold face). Vectors are always considered to be column vectors. The expected value of a random vector is a vector of constants,

notation

$$(6.1.2) \quad \mathcal{E}[\mathbf{x}] = \begin{bmatrix} E[x_1] \\ \vdots \\ E[x_n] \end{bmatrix}$$

For two random variables  $\mathbf{x}$  and  $\mathbf{y}$ , their *covariance* is defined as

$$(6.1.3) \quad \text{cov}[\mathbf{x}, \mathbf{y}] = E[(\mathbf{x} - E[\mathbf{x}])(\mathbf{y} - E[\mathbf{y}])]$$

Computation rules with covariances are

$$(6.1.4) \quad \text{cov}[\mathbf{x}, \mathbf{z}] = \text{cov}[\mathbf{z}, \mathbf{x}] \quad \text{cov}[\mathbf{x}, \mathbf{x}] = \text{var}[\mathbf{x}] \quad \text{cov}[\mathbf{x}, \alpha] = 0$$

$$(6.1.5) \quad \text{cov}[\mathbf{x} + \mathbf{y}, \mathbf{z}] = \text{cov}[\mathbf{x}, \mathbf{z}] + \text{cov}[\mathbf{y}, \mathbf{z}] \quad \text{cov}[\alpha \mathbf{x}, \mathbf{y}] = \alpha \text{cov}[\mathbf{x}, \mathbf{y}]$$

PROBLEM 114. 3 points Using definition (6.1.3) prove the following formula.

$$(6.1.6) \quad \text{cov}[\mathbf{x}, \mathbf{y}] = E[\mathbf{x}\mathbf{y}] - E[\mathbf{x}] E[\mathbf{y}].$$

Write it down carefully, you will lose points for unbalanced or missing paranthes and brackets.

ANSWER. Here it is side by side with and without the notation  $E[\mathbf{x}] = \mu$  and  $E[\mathbf{y}] = \nu$ :

$$(6.1.7) \quad \begin{aligned} \text{cov}[\mathbf{x}, \mathbf{y}] &= E[(\mathbf{x} - E[\mathbf{x}])(\mathbf{y} - E[\mathbf{y}])] & \text{cov}[\mathbf{x}, \mathbf{y}] &= E[(\mathbf{x} - \mu)(\mathbf{y} - \nu)] \\ &= E[\mathbf{x}\mathbf{y} - \mathbf{x} E[\mathbf{y}] - E[\mathbf{x}]\mathbf{y} + E[\mathbf{x}] E[\mathbf{y}]] & &= E[\mathbf{x}\mathbf{y} - \mathbf{x}\nu - \mu\mathbf{y} + \mu\nu] \\ &= E[\mathbf{x}\mathbf{y}] - E[\mathbf{x}] E[\mathbf{y}] - E[\mathbf{x}]\mathbf{y} + E[\mathbf{x}] E[\mathbf{y}] & &= E[\mathbf{x}\mathbf{y}] - \mu\nu - \mu\nu + \mu\nu \\ &= E[\mathbf{x}\mathbf{y}] - E[\mathbf{x}] E[\mathbf{y}]. & &= E[\mathbf{x}\mathbf{y}] - \mu\nu. \end{aligned}$$

PROBLEM 115. 1 point Using (6.1.6) prove the five computation rules with variances (6.1.4) and (6.1.5).

PROBLEM 116. Using the computation rules with covariances, show that

$$(6.1.8) \quad \text{var}[\mathbf{x} + \mathbf{y}] = \text{var}[\mathbf{x}] + 2 \text{cov}[\mathbf{x}, \mathbf{y}] + \text{var}[\mathbf{y}].$$

If one deals with random vectors, the expected value becomes a vector, and the variance becomes a matrix, which is called *dispersion matrix* or *variance-covariance matrix* or simply *covariance matrix*. We will write it  $\mathcal{V}[\mathbf{x}]$ . Its formal definition is

$$(6.1.9) \quad \mathcal{V}[\mathbf{x}] = \mathcal{E}[(\mathbf{x} - \mathcal{E}[\mathbf{x}])(\mathbf{x} - \mathcal{E}[\mathbf{x}])^T],$$

but we can look at it simply as the matrix of all variances and covariances, example

$$(6.1.10) \quad \mathcal{V} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \text{var}[\mathbf{x}] & \text{cov}[\mathbf{x}, \mathbf{y}] \\ \text{cov}[\mathbf{y}, \mathbf{x}] & \text{var}[\mathbf{y}] \end{bmatrix}.$$

An important computation rule for the covariance matrix is

$$(6.1.11) \quad \mathcal{V}[\mathbf{x}] = \mathbf{\Psi} \Rightarrow \mathcal{V}[\mathbf{Ax}] = \mathbf{A}\mathbf{\Psi}\mathbf{A}^\top.$$

PROBLEM 117. 4 points Let  $\mathbf{x} = \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix}$  be a vector consisting of two random variables, with covariance matrix  $\mathcal{V}[\mathbf{x}] = \mathbf{\Psi}$ , and let  $\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$  be an arbitrary  $2 \times 2$  matrix. Prove that

$$(6.1.12) \quad \mathcal{V}[\mathbf{Ax}] = \mathbf{A}\mathbf{\Psi}\mathbf{A}^\top.$$

Hint: You need to multiply matrices, and to use the following computation rules for covariances:

$$(6.1.13) \quad \begin{aligned} \mathcal{Cov}[\mathbf{x} + \mathbf{y}, \mathbf{z}] &= \mathcal{Cov}[\mathbf{x}, \mathbf{z}] + \mathcal{Cov}[\mathbf{y}, \mathbf{z}] & \mathcal{Cov}[\alpha\mathbf{x}, \mathbf{y}] &= \alpha \mathcal{Cov}[\mathbf{x}, \mathbf{y}] & \mathcal{Cov}[\mathbf{x}, \mathbf{x}] &= \text{var}[\mathbf{x}]. \end{aligned}$$

ANSWER.  $\mathcal{V}[\mathbf{Ax}] =$

$$\mathcal{V}\left[\begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} \mathbf{y} \\ \mathbf{z} \end{pmatrix}\right] = \mathcal{V}\left[\begin{pmatrix} a\mathbf{y} + b\mathbf{z} \\ c\mathbf{y} + d\mathbf{z} \end{pmatrix}\right] = \begin{bmatrix} \text{var}[a\mathbf{y} + b\mathbf{z}] & \text{cov}[a\mathbf{y} + b\mathbf{z}, c\mathbf{y} + d\mathbf{z}] \\ \text{cov}[c\mathbf{y} + d\mathbf{z}, a\mathbf{y} + b\mathbf{z}] & \text{var}[c\mathbf{y} + d\mathbf{z}] \end{bmatrix}$$

On the other hand,  $\mathbf{A}\mathbf{\Psi}\mathbf{A}^\top =$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} \text{var}[\mathbf{y}] & \text{cov}[\mathbf{y}, \mathbf{z}] \\ \text{cov}[\mathbf{y}, \mathbf{z}] & \text{var}[\mathbf{z}] \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix} = \begin{bmatrix} a \text{var}[\mathbf{y}] + b \text{cov}[\mathbf{y}, \mathbf{z}] & a \text{cov}[\mathbf{y}, \mathbf{z}] + b \text{var}[\mathbf{z}] \\ c \text{var}[\mathbf{y}] + d \text{cov}[\mathbf{y}, \mathbf{z}] & c \text{cov}[\mathbf{y}, \mathbf{z}] + d \text{var}[\mathbf{z}] \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix}$$

Multiply out and show that it is the same thing.  $\square$

Since the variances are nonnegative, one can see from equation (6.1.11) that covariance matrices are *nonnegative definite* (which is in econometrics is often also called *positive semidefinite*). By definition, a symmetric matrix  $\mathbf{\Sigma}$  is nonnegative definite if for all vectors  $\mathbf{a}$  follows  $\mathbf{a}^\top \mathbf{\Sigma} \mathbf{a} \geq 0$ . It is *positive definite* if it is nonnegative definite, and  $\mathbf{a}^\top \mathbf{\Sigma} \mathbf{a} = 0$  holds only if  $\mathbf{a} = \mathbf{o}$ .

PROBLEM 118. 1 point A symmetric matrix  $\mathbf{\Omega}$  is nonnegative definite if and only if  $\mathbf{a}^\top \mathbf{\Omega} \mathbf{a} \geq 0$  for every vector  $\mathbf{a}$ . Using this criterion, show that if  $\mathbf{\Sigma}$  is symmetric and nonnegative definite, and if  $\mathbf{R}$  is an arbitrary matrix, then  $\mathbf{R}^\top \mathbf{\Sigma} \mathbf{R}$  is also nonnegative definite.

One can also define a covariance matrix between *different* vectors,  $\mathcal{C}[\mathbf{x}, \mathbf{y}]$ ; its  $i, j$  element is  $\text{cov}[x_i, y_j]$ .

The *correlation coefficient* of two scalar random variables is defined as

$$(6.1.14) \quad \text{corr}[x, y] = \frac{\text{cov}[x, y]}{\sqrt{\text{var}[x] \text{var}[y]}}.$$

The advantage of the correlation coefficient over the covariance is that it is always between  $-1$  and  $+1$ . This follows from the Cauchy-Schwartz inequality

$$(6.1.15) \quad (\text{cov}[x, y])^2 \leq \text{var}[x] \text{var}[y].$$

PROBLEM 119. 4 points Given two random variables  $\mathbf{y}$  and  $\mathbf{z}$  with  $\text{var}[\mathbf{y}] \neq 0$  compute that constant  $a$  for which  $\text{var}[a\mathbf{y} - \mathbf{z}]$  is the minimum. Then derive Cauchy-Schwartz inequality from the fact that the minimum variance is nonnegative.

ANSWER.

$$(6.1.16) \quad \text{var}[a\mathbf{y} - \mathbf{z}] = a^2 \text{var}[\mathbf{y}] - 2a \text{cov}[\mathbf{y}, \mathbf{z}] + \text{var}[\mathbf{z}]$$

$$(6.1.17) \quad \text{First order condition: } 0 = 2a \text{var}[\mathbf{y}] - 2 \text{cov}[\mathbf{y}, \mathbf{z}]$$

Therefore the minimum value is  $a^* = \text{cov}[\mathbf{y}, \mathbf{z}] / \text{var}[\mathbf{y}]$ , for which the cross product term is  $-2$  times the first item:

$$(6.1.18) \quad 0 \leq \text{var}[a^*\mathbf{y} - \mathbf{z}] = \frac{(\text{cov}[\mathbf{y}, \mathbf{z}])^2}{\text{var}[\mathbf{y}]} - \frac{2(\text{cov}[\mathbf{y}, \mathbf{z}])^2}{\text{var}[\mathbf{y}]} + \text{var}[\mathbf{z}]$$

$$(6.1.19) \quad 0 \leq -(\text{cov}[\mathbf{y}, \mathbf{z}])^2 + \text{var}[\mathbf{y}] \text{var}[\mathbf{z}].$$

This proves (6.1.15) for the case  $\text{var}[\mathbf{y}] \neq 0$ . If  $\text{var}[\mathbf{y}] = 0$ , then  $\mathbf{y}$  is a constant, therefore  $\text{cov}[\mathbf{y}, \mathbf{z}] = 0$  and (6.1.15) holds trivially.

## 6.2. Marginal Probability Laws

The *marginal* probability distribution of  $\mathbf{x}$  (or  $\mathbf{y}$ ) is simply the probability distribution of  $\mathbf{x}$  (or  $\mathbf{y}$ ). The word “marginal” merely indicates that it is derived from the joint probability distribution of  $\mathbf{x}$  and  $\mathbf{y}$ .

If the probability distribution is characterized by a probability mass function we can compute the marginal probability mass functions by writing down the joint probability mass function in a rectangular scheme and summing up the rows and columns:

$$(6.2.1) \quad p_{\mathbf{x}}(x) = \sum_{\mathbf{y}: p(x, \mathbf{y}) \neq 0} p_{\mathbf{x}, \mathbf{y}}(x, \mathbf{y}).$$

For density functions, the following argument can be given:

$$(6.2.2) \quad \Pr[\mathbf{x} \in dV_x] = \Pr\left[\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \in dV_x \times \mathbb{R}\right].$$

By the definition of a product set:  $\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in A \times B \Leftrightarrow \mathbf{x} \in A \text{ and } \mathbf{y} \in B$ . Split  $\mathbb{R}$  into many small disjoint intervals,  $\mathbb{R} = \bigcup_i dV_{y_i}$ , then

$$(6.2.3) \quad \Pr[\mathbf{x} \in dV_x] = \sum_i \Pr\left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in dV_x \times dV_{y_i}\right]$$

$$(6.2.4) \quad = \sum_i f_{\mathbf{x},\mathbf{y}}(x, y_i) |dV_x| |dV_{y_i}|$$

$$(6.2.5) \quad = |dV_x| \sum_i f_{\mathbf{x},\mathbf{y}}(x, y_i) |dV_{y_i}|.$$

Therefore  $\sum_i f_{\mathbf{x},\mathbf{y}}(x, y) |dV_{y_i}|$  is the density function we are looking for. Now the  $|dV_{y_i}|$  are usually written as  $dy$ , and the sum is usually written as an integral (i.e., an infinite sum each summand of which is infinitesimal), therefore we get

$$(6.2.6) \quad f_{\mathbf{x}}(x) = \int_{y=-\infty}^{y=+\infty} f_{\mathbf{x},\mathbf{y}}(x, y) dy.$$

In other words, one has to “integrate out” the variable which one is not interested in.

### 6.3. Conditional Probability Distribution and Conditional Mean

The conditional probability distribution of  $\mathbf{y}$  given  $\mathbf{x}=x$  is the probability distribution of  $\mathbf{y}$  if we count only those experiments in which the outcome of  $\mathbf{x}$  is  $x$ . If the distribution is defined by a probability mass function, then this is no problem:

$$(6.3.1) \quad p_{\mathbf{y}|\mathbf{x}}(y, x) = \Pr[\mathbf{y}=y|\mathbf{x}=x] = \frac{\Pr[\mathbf{y}=y \text{ and } \mathbf{x}=x]}{\Pr[\mathbf{x}=x]} = \frac{p_{\mathbf{x},\mathbf{y}}(x, y)}{p_{\mathbf{x}}(x)}.$$

For a density function there is the problem that  $\Pr[\mathbf{x}=x] = 0$ , i.e., the conditional probability is strictly speaking not defined. Therefore take an infinitesimal volume element  $dV_x$  located at  $x$  and condition on  $\mathbf{x} \in dV_x$ :

$$(6.3.2) \quad \Pr[\mathbf{y} \in dV_y | \mathbf{x} \in dV_x] = \frac{\Pr[\mathbf{y} \in dV_y \text{ and } \mathbf{x} \in dV_x]}{\Pr[\mathbf{x} \in dV_x]}$$

$$(6.3.3) \quad = \frac{f_{\mathbf{x},\mathbf{y}}(x, y) |dV_x| |dV_y|}{f_{\mathbf{x}}(x) |dV_x|}$$

$$(6.3.4) \quad = \frac{f_{\mathbf{x},\mathbf{y}}(x, y)}{f_{\mathbf{x}}(x)} |dV_y|.$$

This no longer depends on  $dV_x$ , only on its location  $x$ . The conditional density therefore

$$(6.3.5) \quad f_{\mathbf{y}|\mathbf{x}}(y, x) = \frac{f_{\mathbf{x},\mathbf{y}}(x, y)}{f_{\mathbf{x}}(x)}.$$

As  $y$  varies, the conditional density is proportional to the joint density function, but for every given value of  $x$  the joint density is multiplied by an appropriate factor that its integral with respect to  $y$  is 1. From (6.3.5) follows also that the joint density function is the product of the conditional times the marginal density functions.

**PROBLEM 120.** 2 points *The conditional density is the joint divided by the marginal:*

$$(6.3.6) \quad f_{\mathbf{y}|\mathbf{x}}(y, x) = \frac{f_{\mathbf{x},\mathbf{y}}(x, y)}{f_{\mathbf{x}}(x)}.$$

*Show that this density integrates out to 1.*

**ANSWER.** The conditional is a density in  $y$  with  $x$  as parameter. Therefore its integral with respect to  $y$  must be = 1. Indeed,

$$(6.3.7) \quad \int_{y=-\infty}^{y=+\infty} f_{\mathbf{y}|\mathbf{x}=x}(y, x) dy = \frac{\int_{y=-\infty}^{y=+\infty} f_{\mathbf{x},\mathbf{y}}(x, y) dy}{f_{\mathbf{x}}(x)} = \frac{f_{\mathbf{x}}(x)}{f_{\mathbf{x}}(x)} = 1$$

because of the formula for the marginal:

$$(6.3.8) \quad f_{\mathbf{x}}(x) = \int_{y=-\infty}^{y=+\infty} f_{\mathbf{x},\mathbf{y}}(x, y) dy$$

You see that formula (6.3.6) divides the joint density exactly by the right number which makes integral equal to 1.

**PROBLEM 121.** [BD77, example 1.1.4 on p. 7].  $x$  and  $y$  are two independent random variables uniformly distributed over  $[0, 1]$ . Define  $\mathbf{u} = \min(x, y)$  and  $\mathbf{v} = \max(x, y)$ .

• a. *Draw in the  $x, y$  plane the event  $\{\max(x, y) \leq 0.5 \text{ and } \min(x, y) > 0.4\}$  and compute its probability.*

**ANSWER.** The event is the square between 0.4 and 0.5, and its probability is 0.01.

• b. *Compute the probability of the event  $\{\max(x, y) \leq 0.5 \text{ and } \min(x, y) \leq 0.4\}$ .*

**ANSWER.** It is  $\Pr[\max(x, y) \leq 0.5] - \Pr[\max(x, y) \leq 0.5 \text{ and } \min(x, y) > 0.4]$ , i.e., the area of the square from 0 to 0.5 minus the square we just had, i.e., 0.24.

• c. *Compute  $\Pr[\max(x, y) \leq 0.5 | \min(x, y) \leq 0.4]$ .*

ANSWER.

$$(6.3.9) \quad \frac{\Pr[\max(x, y) \leq 0.5 \text{ and } \min(x, y) \leq 0.4]}{\Pr[\min(x, y) \leq 0.4]} = \frac{0.24}{1 - 0.36} = \frac{0.24}{0.64} = \frac{3}{8}.$$

□

- d. Compute the joint cumulative distribution function of  $\mathbf{u}$  and  $\mathbf{v}$ .

ANSWER. One good way is to do it geometrically: for arbitrary  $0 \leq u, v \leq 1$  draw the area  $\{\mathbf{u} \leq u \text{ and } \mathbf{v} \leq v\}$  and then derive its size. If  $u \leq v$  then  $\Pr[\mathbf{u} \leq u \text{ and } \mathbf{v} \leq v] = \Pr[\mathbf{v} \leq v] - \Pr[\mathbf{u} \leq u \text{ and } \mathbf{v} > v] = v^2 - (v - u)^2 = 2uv - u^2$ . If  $u \geq v$  then  $\Pr[\mathbf{u} \leq u \text{ and } \mathbf{v} \leq v] = \Pr[\mathbf{v} \leq v] = v^2$ .

□

- e. Compute the joint density function of  $\mathbf{u}$  and  $\mathbf{v}$ . Note: this joint density is discontinuous. The values at the breakpoints themselves do not matter, but it is very important to give the limits within this is a nontrivial function and where it is zero.

ANSWER. One can see from the way the cumulative distribution function was constructed that the density function must be

$$(6.3.10) \quad f_{\mathbf{u}, \mathbf{v}}(u, v) = \begin{cases} 2 & \text{if } 0 \leq u \leq v \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

I.e., it is uniform in the above-diagonal part of the square. This is also what one gets from differentiating  $2vu - u^2$  once with respect to  $u$  and once with respect to  $v$ .

□

- f. Compute the marginal density function of  $\mathbf{u}$ .

ANSWER. Integrate  $v$  out: the marginal density of  $\mathbf{u}$  is

$$(6.3.11) \quad f_{\mathbf{u}}(u) = \int_{v=u}^1 2dv = 2v \Big|_u^1 = 2 - 2u \quad \text{if } 0 \leq u \leq 1, \quad \text{and } 0 \text{ otherwise.}$$

□

- g. Compute the conditional density of  $\mathbf{v}$  given  $\mathbf{u} = u$ .

ANSWER. Conditional density is easy to get too; it is the joint divided by the marginal, i.e., it is uniform:

$$(6.3.12) \quad f_{\mathbf{v}|\mathbf{u}=u}(v) = \begin{cases} \frac{1}{1-u} & \text{for } 0 \leq u \leq v \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

□

### 6.4. The Multinomial Distribution

Assume you have an experiment with  $r$  different possible outcomes, with outcome  $i$  having probability  $p_i$  ( $i = 1, \dots, r$ ). You are repeating the experiment  $n$  different times, and you count how many times the  $i$ th outcome occurred. Therefore you get

a random vector with  $r$  different components  $x_i$ , indicating how often the  $i$ th event occurred. The probability to get the frequencies  $x_1, \dots, x_r$  is

$$(6.4.1) \quad \Pr[\mathbf{x}_1 = x_1, \dots, \mathbf{x}_r = x_r] = \frac{m!}{x_1! \cdots x_r!} p_1^{x_1} p_2^{x_2} \cdots p_r^{x_r}$$

This can be explained as follows: The probability that the first  $x_1$  experiments yield outcome 1, the next  $x_2$  outcome 2, etc., is  $p_1^{x_1} p_2^{x_2} \cdots p_r^{x_r}$ . Now every other sequence of experiments which yields the same number of outcomes of the different categories is simply a permutation of this. But multiplying this probability by  $n!$  may count certain sequences of outcomes more than once. Therefore we have to divide by the number of permutations of the whole  $n$  element set which yield the same original sequence. This is  $x_1! \cdots x_r!$ , because this must be a permutation which permutes the first  $x_1$  elements amongst themselves, etc. Therefore the relevant count of permutations is  $\frac{n!}{x_1! \cdots x_r!}$ .

PROBLEM 122. You have an experiment with  $r$  different outcomes, the  $i$ th outcome occurring with probability  $p_i$ . You make  $n$  independent trials, and the  $i$ th outcome occurred  $x_i$  times. The joint distribution of the  $x_1, \dots, x_r$  is called a multinomial distribution with parameters  $n$  and  $p_1, \dots, p_r$ .

- a. 3 points Prove that their mean vector and covariance matrix are

$$(6.4.2) \quad \boldsymbol{\mu} = \mathcal{E} \begin{bmatrix} x_1 \\ \vdots \\ x_r \end{bmatrix} = n \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_r \end{bmatrix} \quad \text{and} \quad \boldsymbol{\Psi} = \mathcal{V} \begin{bmatrix} x_1 \\ \vdots \\ x_r \end{bmatrix} = n \begin{bmatrix} p_1 - p_1^2 & -p_1 p_2 & \cdots & -p_1 p_r \\ -p_2 p_1 & p_2 - p_2^2 & \cdots & -p_2 p_r \\ \vdots & \vdots & \ddots & \vdots \\ -p_r p_1 & -p_r p_2 & \cdots & p_r - p_r^2 \end{bmatrix}$$

Hint: use the fact that the multinomial distribution with parameters  $n$  and  $p_1, \dots, p_r$  is the independent sum of  $n$  multinomial distributions with parameters 1 and  $p_1, \dots, p_r$ .

ANSWER. In one trial,  $x_i^2 = x_i$ , from which follows the formula for the variance, and for  $i \neq j$ ,  $x_i x_j = 0$ , since only one of them can occur. Therefore  $\text{cov}[x_i, x_j] = 0 - \mathcal{E}[x_i] \mathcal{E}[x_j]$ . For several independent trials, just add this.

- b. 1 point How can you show that this covariance matrix is singular?

ANSWER. Since  $x_1 + \cdots + x_r = n$  with zero variance, we should expect

$$(6.4.3) \quad n \begin{bmatrix} p_1 - p_1^2 & -p_1 p_2 & \cdots & -p_1 p_r \\ -p_2 p_1 & p_2 - p_2^2 & \cdots & -p_2 p_r \\ \vdots & \vdots & \ddots & \vdots \\ -p_r p_1 & -p_r p_2 & \cdots & p_r - p_r^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

### 6.5. Independent Random Vectors

The same definition of independence, which we already encountered with scalar random variables, also applies to vector random variables: the vector random variables  $\mathbf{x} : U \rightarrow \mathbb{R}^m$  and  $\mathbf{y} : U \rightarrow \mathbb{R}^n$  are called independent if all events that can be defined in terms of  $\mathbf{x}$  are independent of all events that can be defined in terms of  $\mathbf{y}$ , i.e., all events of the form  $\{\mathbf{x}(\omega) \in C\}$  are independent of all events of the form  $\{\mathbf{y}(\omega) \in D\}$  with arbitrary (measurable) subsets  $C \subset \mathbb{R}^m$  and  $D \subset \mathbb{R}^n$ .

For this it is sufficient that for all  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{y} \in \mathbb{R}^n$ , the event  $\{\mathbf{x} \leq \mathbf{x}\}$  is independent of the event  $\{\mathbf{y} \leq \mathbf{y}\}$ , i.e., that the joint cumulative distribution function is the product of the marginal ones.

Since the joint cumulative distribution function of independent variables is equal to the product of the univariate cumulative distribution functions, the same is true for the joint density function and the joint probability mass function.

Only under this strong definition of independence is it true that any functions of independent random variables are independent.

**PROBLEM 123.** *4 points Prove that, if  $x$  and  $y$  are independent, then  $E[xy] = E[x]E[y]$  and therefore  $\text{cov}[x, y] = 0$ . (You may assume  $x$  and  $y$  have density functions). Give a counterexample where the covariance is zero but the variables are nevertheless dependent.*

**ANSWER.** Just use that the joint density function is the product of the marginals. It can also be done as follows:  $E[xy] = E[E[xy|x]] = E[x E[y|x]] =$  now independence is needed  $= E[x E[y]] = E[x]E[y]$ . A counterexample is given in Problem 139.  $\square$

**PROBLEM 124.** *3 points Prove the following: If the scalar random variables  $x$  and  $y$  are indicator variables (i.e., if each of them can only assume the values 0 and 1), and if  $\text{cov}[x, y] = 0$ , then  $x$  and  $y$  are independent. (I.e., in this respect indicator variables have similar properties as jointly normal random variables.)*

**ANSWER.** Define the events  $A = \{\omega \in U : x(\omega) = 1\}$  and  $B = \{\omega \in U : y(\omega) = 1\}$ , i.e.,  $x = i_A$  (the indicator variable of the event  $A$ ) and  $y = i_B$ . Then  $xy = i_{A \cap B}$ . If  $\text{cov}[x, y] = E[xy] - E[x]E[y] = \Pr[A \cap B] - \Pr[A]\Pr[B] = 0$ , then  $A$  and  $B$  are independent.  $\square$

**PROBLEM 125.** *If the vector random variables  $\mathbf{x}$  and  $\mathbf{y}$  have the property that  $x_i$  is independent of every  $y_j$  for all  $i$  and  $j$ , does that make  $\mathbf{x}$  and  $\mathbf{y}$  independent random vectors? Interestingly, the answer is no. Give a counterexample that this fact does not even hold for indicator variables. I.e., construct two random vectors  $\mathbf{x}$  and  $\mathbf{y}$ , consisting of indicator variables, with the property that each component of  $\mathbf{x}$  is independent of each component of  $\mathbf{y}$ , but  $\mathbf{x}$  and  $\mathbf{y}$  are not independent as vector random variables. Hint: Such an example can be constructed in the simplest possible case that  $\mathbf{x}$  has two components and  $\mathbf{y}$  has one component; i.e., you merely have to find three indicator variables  $x_1, x_2$ , and  $y$  with the property that  $x_1$  is independent*

of  $y$ , and  $x_2$  is independent of  $y$ , but the vector  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  is not independent of  $y$ . In these three variables, you should use three events which are pairwise independent but not mutually independent.

**ANSWER.** Go back to throwing a coin twice independently and define  $A = \{HH, HT\}$ ;  $B = \{TH, HH\}$ , and  $C = \{HH, TT\}$ , and  $x_1 = I_A$ ,  $x_2 = I_B$ , and  $y = I_C$ . They are pairwise independent but  $A \cap B \cap C = A \cap B$ , i.e.,  $x_1 x_2 y = x_1 x_2$ , therefore  $E[x_1 x_2 y] \neq E[x_1 x_2] E[y]$  therefore they are not independent.

**PROBLEM 126.** *4 points Prove that, if  $x$  and  $y$  are independent, then  $\text{var}[xy] = (E[x])^2 \text{var}[y] + (E[y])^2 \text{var}[x] + \text{var}[x] \text{var}[y]$ .*

**ANSWER.** Start with result and replace all occurrences of  $\text{var}[z]$  with  $E[z^2] - E[z]^2$ , then multiply out:  $E[x]^2(E[y^2] - E[y]^2) + E[y]^2(E[x^2] - E[x]^2) + (E[x^2] - E[x]^2)(E[y^2] - E[y]^2) = E[x^2]E[y^2] - E[x]^2E[y]^2 = E[(xy)^2] - E[xy]^2$ .

### 6.6. Conditional Expectation and Variance

The *conditional expectation* of  $y$  is the expected value of  $y$  under the conditional density. If joint densities exist, it follows

$$(6.6.1) \quad E[y|x=x] = \frac{\int y f_{x,y}(x, y) dy}{f_x(x)} =: g(x).$$

This is not a random variable but a constant which depends on  $x$ , i.e., a function of  $x$ , which is called here  $g(x)$ . But often one uses the term  $E[y|x]$  without specifying  $x$ . This is, by definition, the random variable  $g(x)$  which one gets by plugging  $x$  in  $g$ ; it assigns to every outcome  $\omega \in U$  the conditional expectation of  $y$  given  $x=x(\omega)$ .

Since  $E[y|x]$  is a random variable, it is possible to take its expected value. The law of iterated expectations is extremely important here. It says that you will get the same result as if you had taken the expected value of  $y$ :

$$(6.6.2) \quad E[E[y|x]] = E[y].$$

Proof (for the case that the densities exist):

$$(6.6.3) \quad \begin{aligned} E[E[y|x]] &= E[g(x)] = \int \frac{\int y f_{x,y}(x, y) dy}{f_x(x)} f_x(x) dx \\ &= \int \int y f_{x,y}(x, y) dy dx = E[y]. \end{aligned}$$

**PROBLEM 127.** *Let  $x$  and  $y$  be two jointly distributed variables. For every fixed value  $x$ ,  $\text{var}[y|x = x]$  is the variance of  $y$  under the conditional distribution, and  $\text{var}[y|x]$  is this variance as a random variable, namely, as a function of  $x$ .*



- a. 1 point Prove that

$$(6.6.4) \quad \text{var}[y|x] = E[y^2|x] - (E[y|x])^2.$$

This is a very simple proof. Explain exactly what, if anything, needs to be done to prove it.

ANSWER. For every fixed value  $x$ , it is an instance of the law

$$(6.6.5) \quad \text{var}[y] = E[y^2] - (E[y])^2$$

applied to the conditional density given  $x = x$ . And since it is true for every fixed  $x$ , it is also true after plugging in the random variable  $x$ .  $\square$

- b. 3 points Prove that

$$(6.6.6) \quad \text{var}[y] = \text{var}[E[y|x]] + E[\text{var}[y|x]],$$

i.e., the variance consists of two components: the variance of the conditional mean and the mean of the conditional variances. This decomposition of the variance is given e.g. in [Rao73, p. 97] or [Ame94, theorem 4.4.2 on p. 78].

ANSWER. The first term on the rhs is  $E[(E[y|x])^2] - (E[E[y|x]])^2$ , and the second term, due to (6.6.4), becomes  $E[E[y^2|x] - E[(E[y|x])^2]]$ . If one adds, the two  $E[(E[y|x])^2]$  cancel out, and the other two terms can be simplified by the law of iterated expectations to give  $E[y^2] - (E[y])^2$ .  $\square$

- c. 2 points [Coo98, p. 23] The conditional expected value is sometimes called the population regression function. In graphical data analysis, the sample equivalent of the variance ratio

$$(6.6.7) \quad \frac{E[\text{var}[y|x]]}{\text{var}[E[y|x]]}$$

can be used to determine whether the regression function  $E[y|x]$  appears to be visually well-determined or not. Does a small or a big variance ratio indicate a well-determined regression function?

ANSWER. For a well-determined regression function the variance ratio should be small. [Coo98, p. 23] writes: “This ratio is reminiscent of a one-way analysis of variance, with the numerator representing the average within group (slice) variance, and the denominator representing the variance between group (slice) means.”  $\square$

Now some general questions:

PROBLEM 128. The figure on page 102 shows 250 independent observations of the random vector  $\begin{bmatrix} x \\ y \end{bmatrix}$ .

- a. 2 points Draw in by hand the approximate location of  $\mathcal{E}[\begin{bmatrix} x \\ y \end{bmatrix}]$  and the graph of  $E[y|x]$ . Draw into the second diagram the approximate marginal density of  $x$ .

- b. 2 points Is there a law that the graph of the conditional expectation  $E[y|x]$  always goes through the point  $\mathcal{E}[\begin{bmatrix} x \\ y \end{bmatrix}]$ —for arbitrary probability distributions for which these expectations exist, or perhaps for an important special case? Indicate how this could be proved or otherwise give (maybe geometrically) a simple counterexample.

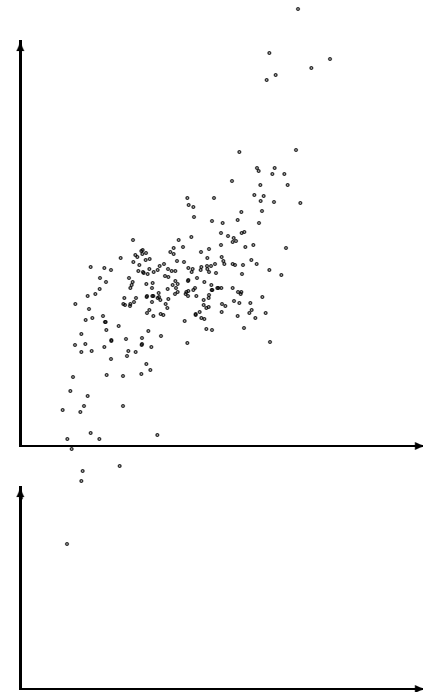
ANSWER. This is *not* the law of iterated expectations. It is true for jointly normal variables, but not in general. It is also true if  $x$  and  $y$  are independent; then the graph of  $E[y|x]$  is a horizontal line at the height of the unconditional expectation  $E[y]$ . A distribution with U-shaped unconditional density has the unconditional mean in the center of the U, i.e., here the unconditional mean does not lie on the curve drawn out by the conditional mean.

- c. 2 points Do you have any ideas how the strange-looking cluster of points in the figure on page 102 was generated?

PROBLEM 129. 2 points Given two independent random variables  $x$  and  $y$  with marginal density functions  $f_x(x)$  and  $g_y(y)$ . Write down their joint, marginal, and conditional densities.

ANSWER. Joint density:  $f_{x,y}(x, y) = f_x(x)g_y(y)$ .

Marginal density of  $x$  is  $\int_{-\infty}^{\infty} f_x(x)g_y(y) dy = f_x(x) \int_{-\infty}^{\infty} g_y(y) dy = f_x(x)$ , and that of  $y$  is  $\int_{-\infty}^{\infty} f_x(x)g_y(y) dx = g_y(y) \int_{-\infty}^{\infty} f_x(x) dx = g_y(y)$ . The text of the question should have been: “Given two independent random variables  $x$  and  $y$  with marginal density functions  $f_x(x)$  and  $g_y(y)$ ”; by just calling them “density functions” without specifying “marginal” it committed the error of de-totalization, i.e., it treated elements



a totality, i.e., of an ensemble in which each depends on everything else, as if they could be defined independently of each other.

Conditional density functions:  $f_{x|y=y}(x; y) = f_x(x)$  (i.e., it does not depend on  $y$ ); and  $g_{y|x=x}(y; x) = g_y(y)$ . You can see this by dividing the joint by the marginal.  $\square$

### 6.7. Expected Values as Predictors

Expected values and conditional expected values have optimal properties as predictors.

**PROBLEM 130.** *3 points* What is the best predictor of a random variable  $\mathbf{y}$  by a constant  $a$ , if the loss function is the “mean squared error” (MSE)  $E[(\mathbf{y} - a)^2]$ ?

ANSWER. Write  $E[\mathbf{y}] = \mu$ ; then

$$(6.7.1) \quad \begin{aligned} (\mathbf{y} - a)^2 &= ((\mathbf{y} - \mu) - (a - \mu))^2 \\ &= (\mathbf{y} - \mu)^2 - 2(\mathbf{y} - \mu)(a - \mu) + (a - \mu)^2; \end{aligned}$$

$$\text{therefore } E[(\mathbf{y} - a)^2] = E[(\mathbf{y} - \mu)^2] - 0 + (a - \mu)^2$$

This is minimized by  $a = \mu$ .  $\square$

The expected value of  $\mathbf{y}$  is therefore that constant which, as predictor of  $\mathbf{y}$ , has smallest MSE.

What if we want to predict  $\mathbf{y}$  not by a constant but by a function of the random vector  $\mathbf{x}$ , call it  $h(\mathbf{x})$ ?

**PROBLEM 131.** *2 points* Assume the vector  $\mathbf{x} = [x_1, \dots, x_j]^\top$  and the scalar  $\mathbf{y}$  are jointly distributed random variables, and assume conditional means exist.  $\mathbf{x}$  is observed, but  $\mathbf{y}$  is not observed. The joint distribution of  $\mathbf{x}$  and  $\mathbf{y}$  is known. Show that the conditional expectation  $E[\mathbf{y}|\mathbf{x}]$  is the minimum MSE predictor of  $\mathbf{y}$  given  $\mathbf{x}$ , i.e., show that for any other function of  $\mathbf{x}$ , call it  $h(\mathbf{x})$ , the following inequality holds:

$$(6.7.2) \quad E[(\mathbf{y} - h(\mathbf{x}))^2] \geq E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])^2].$$

For this proof and the proofs required in Problems 132 and 133, you may use (1) the theorem of iterated expectations  $E[E[\mathbf{y}|\mathbf{x}]] = E[\mathbf{y}]$ , (2) the additivity  $E[g(\mathbf{y}) + h(\mathbf{y})|\mathbf{x}] = E[g(\mathbf{y})|\mathbf{x}] + E[h(\mathbf{y})|\mathbf{x}]$ , and (3) the fact that  $E[g(\mathbf{x})h(\mathbf{y})|\mathbf{x}] = g(\mathbf{x})E[h(\mathbf{y})|\mathbf{x}]$ . Be very specific about which rules you are applying at every step. You must show that you understand what you are writing down.

ANSWER.

$$(6.7.3) \quad \begin{aligned} E[(\mathbf{y} - h(\mathbf{x}))^2] &= E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}] - (h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}]))^2] \\ &= E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])^2] - 2E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])(h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])] + E[(h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])^2]. \end{aligned}$$

Here the cross product term  $E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])(h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])]$  is zero. In order to see this, first use the law of iterated expectations

$$(6.7.4) \quad E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])(h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])] = E[E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])(h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])|\mathbf{x}]]$$

and then look at the inner term, not yet doing the outer expectation:

$$E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])(h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])|\mathbf{x}] = (h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])$$

$$E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])|\mathbf{x}] = (h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])(E[\mathbf{y}|\mathbf{x}] - E[\mathbf{y}|\mathbf{x}]) = (h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}]) \cdot 0 = 0$$

Plugging this into (6.7.4) gives  $E[(\mathbf{y} - E[\mathbf{y}|\mathbf{x}])(h(\mathbf{x}) - E[\mathbf{y}|\mathbf{x}])] = E[0] = 0$ .

This is one of the few clear cut results in probability theory where a best estimator/predictor exists. In this case, however, all parameters of the distribution are known, the only uncertainty comes from the fact that some random variables are unobserved.

**PROBLEM 132.** Assume the vector  $\mathbf{x} = [x_1, \dots, x_j]^\top$  and the scalar  $\mathbf{y}$  are jointly distributed random variables, and assume conditional means exist. Define  $\varepsilon = \mathbf{y} - E[\mathbf{y}|\mathbf{x}]$ .

• a. *5 points* Demonstrate the following identities:

$$(6.7.5) \quad E[\varepsilon|\mathbf{x}] = 0$$

$$(6.7.6) \quad E[\varepsilon] = 0$$

$$(6.7.7) \quad E[x_i \varepsilon|\mathbf{x}] = 0 \quad \text{for all } i, 1 \leq i \leq j$$

$$(6.7.8) \quad E[x_i \varepsilon] = 0 \quad \text{for all } i, 1 \leq i \leq j$$

$$(6.7.9) \quad \text{cov}[x_i, \varepsilon] = 0 \quad \text{for all } i, 1 \leq i \leq j.$$

Interpretation of (6.7.9):  $\varepsilon$  is the error in the best prediction of  $\mathbf{y}$  based on  $\mathbf{x}$ . If the error were correlated with one of the components  $x_i$ , then this correlation could be used to construct a better prediction of  $\mathbf{y}$ .

ANSWER. (6.7.5):  $E[\varepsilon|\mathbf{x}] = E[\mathbf{y}|\mathbf{x}] - E[E[\mathbf{y}|\mathbf{x}]|\mathbf{x}] = 0$  since  $E[\mathbf{y}|\mathbf{x}]$  is a function of  $\mathbf{x}$  and therefore equal to its own expectation conditionally on  $\mathbf{x}$ . (This is *not* the law of iterated expectations; the law that the expected value of a constant is a constant.)

(6.7.6) follows from (6.7.5) (i.e., (6.7.5) is stronger than (6.7.6)): if an expectation is zero conditionally on every possible outcome of  $\mathbf{x}$  then it is zero altogether. In formulas,  $E[\varepsilon] = E[E[\varepsilon|\mathbf{x}]] = E[0] = 0$ . It is also easy to show it in one swoop, without using (6.7.5):  $E[\varepsilon] = E[\mathbf{y} - E[\mathbf{y}|\mathbf{x}]] = 0$ . Either way you need the law of iterated expectations for this.

$$(6.7.7): E[x_i \varepsilon|\mathbf{x}] = x_i E[\varepsilon|\mathbf{x}] = 0.$$

(6.7.8):  $E[x_i \varepsilon] = E[E[x_i \varepsilon|\mathbf{x}]] = E[0] = 0$ ; or in one swoop:  $E[x_i \varepsilon] = E[x_i \mathbf{y} - x_i E[\mathbf{y}|\mathbf{x}]] = E[x_i \mathbf{y}] - E[x_i \mathbf{y}] = E[x_i \mathbf{y}] - E[x_i \mathbf{y}] = 0$ . The following “proof” is not correct:  $E[x_i \varepsilon] = E[x_i] E[\varepsilon] = E[x_i] \cdot 0 = 0$ .  $x_i$  and  $\varepsilon$  are generally not independent, therefore the multiplication rule  $E[x_i \varepsilon] = E[x_i] E[\varepsilon]$  cannot be used. Of course, the following “proof” does not work either:  $E[x_i \varepsilon] = x_i E[\varepsilon] = x_i \cdot 0 = 0$ .  $x_i$  is a random variable and  $E[x_i \varepsilon]$  is a constant; therefore  $E[x_i \varepsilon] = x_i E[\varepsilon]$  cannot hold.

$$(6.7.9): \text{cov}[x_i, \varepsilon] = E[x_i \varepsilon] - E[x_i] E[\varepsilon] = 0 - E[x_i] \cdot 0 = 0.$$

• b. *2 points* This part can only be done after discussing the multivariate normal distribution. If  $\mathbf{x}$  and  $\mathbf{y}$  are jointly normal, show that  $\mathbf{x}$  and  $\varepsilon$  are independent, and

that the variance of  $\varepsilon$  does not depend on  $\mathbf{x}$ . (This is why one can consider it an error term.)

ANSWER. If  $\mathbf{x}$  and  $\mathbf{y}$  are jointly normal, then  $\mathbf{x}$  and  $\varepsilon$  are jointly normal as well, and independence follows from the fact that their covariance is zero. The variance is constant because in the Normal case, the conditional variance is constant, i.e.,  $E[\varepsilon^2] = E[E[\varepsilon^2|\mathbf{x}]] = \text{constant}$  (does not depend on  $\mathbf{x}$ ).  $\square$

PROBLEM 133. 5 points Under the permanent income hypothesis, the assumption is made that consumers' lifetime utility is highest if the same amount is consumed every year. The utility-maximizing level of consumption  $\mathbf{c}$  for a given consumer depends on the actual state of the economy in each of the  $n$  years of the consumer's life  $\mathbf{c} = f(\mathbf{y}_1, \dots, \mathbf{y}_n)$ . Since  $\mathbf{c}$  depends on future states of the economy, which are not known, it is impossible for the consumer to know this optimal  $\mathbf{c}$  in advance; but it is assumed that the function  $f$  and the joint distribution of  $\mathbf{y}_1, \dots, \mathbf{y}_n$  are known to him. Therefore in period  $t$ , when he only knows the values of  $\mathbf{y}_1, \dots, \mathbf{y}_t$ , but not yet the future values, the consumer decides to consume the amount  $\mathbf{c}_t = E[\mathbf{c}|\mathbf{y}_1, \dots, \mathbf{y}_t]$ , which is the best possible prediction of  $\mathbf{c}$  given the information available to him. Show that in this situation,  $\mathbf{c}_{t+1} - \mathbf{c}_t$  is uncorrelated with all  $\mathbf{y}_1, \dots, \mathbf{y}_t$ . This implication of the permanent income hypothesis can be tested empirically, see [Hal78]. Hint: you are allowed to use without proof the following extension of the theorem of iterated expectations:

$$(6.7.10) \quad E[E[\mathbf{x}|\mathbf{y}, \mathbf{z}|\mathbf{y}]] = E[\mathbf{x}|\mathbf{y}].$$

Here is an explanation of (6.7.10):  $E[\mathbf{x}|\mathbf{y}]$  is the best predictor of  $\mathbf{x}$  based on information set  $\mathbf{y}$ .  $E[\mathbf{x}|\mathbf{y}, \mathbf{z}]$  is the best predictor of  $\mathbf{x}$  based on the extended information set consisting of  $\mathbf{y}$  and  $\mathbf{z}$ .  $E[E[\mathbf{x}|\mathbf{y}, \mathbf{z}|\mathbf{y}]]$  is therefore my prediction, based on  $\mathbf{y}$  only, how I will refine my prediction when  $\mathbf{z}$  becomes available as well. Its equality with  $E[\mathbf{x}|\mathbf{y}]$ , i.e., (6.7.10) says therefore that I cannot predict how I will change my mind after better information becomes available.

ANSWER. In (6.7.10) set  $\mathbf{x} = \mathbf{c} = f(\mathbf{y}_1, \dots, \mathbf{y}_t, \mathbf{y}_{t+1}, \dots, \mathbf{y}_n)$ ,  $\mathbf{y} = [\mathbf{y}_1, \dots, \mathbf{y}_t]^\top$ , and  $\mathbf{z} = \mathbf{y}_{t+1}$  to get

$$(6.7.11) \quad E[E[\mathbf{c}|\mathbf{y}_1, \dots, \mathbf{y}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_t]] = E[\mathbf{c}|\mathbf{y}_1, \dots, \mathbf{y}_t].$$

Writing  $\mathbf{c}_t$  for  $E[\mathbf{c}|\mathbf{y}_1, \dots, \mathbf{y}_t]$ , this becomes  $E[\mathbf{c}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_t] = \mathbf{c}_t$ , i.e.,  $\mathbf{c}_t$  is not only the best predictor of  $\mathbf{c}$ , but also that of  $\mathbf{c}_{t+1}$ . The change in consumption  $\mathbf{c}_{t+1} - \mathbf{c}_t$  is therefore the prediction error, which is uncorrelated with the conditioning variables, as shown in Problem 132.  $\square$

PROBLEM 134. 3 points Show that for any two random variables  $x$  and  $y$  whose covariance exists, the following equation holds:

$$(6.7.12) \quad \text{cov}[x, y] = \text{cov}[x, E[y|x]]$$

Note: Since  $E[y|x]$  is the best predictor of  $y$  based on the observation of  $x$ , (6.7.12) can also be written as

$$(6.7.13) \quad \text{cov}[x, (\mathbf{y} - E[y|x])] = 0,$$

i.e.,  $x$  is uncorrelated with the prediction error of the best prediction of  $y$  given  $x$ . (Nothing to prove for this Note.)

ANSWER. Apply (6.1.6) to the righthand side of (6.7.12):

$$(6.7.14) \quad \text{cov}[x, E[y|x]] = E[xE[y|x]] - E[x]E[E[y|x]] = E[E[xy|x]] - E[x]E[y] = E[xy] - E[x]E[y] = \text{cov}[x, y]$$

The tricky part here is to see that  $x E[y|x] = E[xy|x]$ .

PROBLEM 135. Assume  $x$  and  $y$  have a joint density function  $f_{x,y}(x, y)$  which is symmetric about the  $x$ -axis, i.e.,

$$f_{x,y}(x, y) = f_{x,y}(x, -y).$$

Also assume that variances and covariances exist. Show that  $\text{cov}[x, y] = 0$ . Hint: one way to do it is to look at  $E[y|x]$ .

ANSWER. We know that  $\text{cov}[x, y] = \text{cov}[x, E[y|x]]$ . Furthermore, from symmetry follow  $E[y|x] = 0$ . Therefore  $\text{cov}[x, y] = \text{cov}[x, 0] = 0$ . Here is a detailed proof of  $E[y|x] = 0$ :  $E[y|x=x] = \int_{-\infty}^{\infty} y \frac{f_{x,y}(x, y)}{f_x(x)} dy$ . Now substitute  $z = -y$ , then also  $dz = -dy$ , and the boundaries of integration are reversed:

$$(6.7.15) \quad E[y|x=x] = \int_{-\infty}^{\infty} z \frac{f_{x,y}(x, -z)}{f_x(x)} dz = \int_{\infty}^{-\infty} z \frac{f_{x,y}(x, z)}{f_x(x)} dz = -E[y|x=x].$$

One can also prove directly under this presupposition  $\text{cov}[x, y] = \text{cov}[x, -y]$  and therefore it must be zero.

PROBLEM 136. [Wit85, footnote on p. 241] Let  $\mathbf{p}$  be the logarithm of the price level,  $\mathbf{m}$  the logarithm of the money supply, and  $x$  a variable representing real income. Assume that  $\mathbf{p}$  and  $x$  are independent normal random variables with expected values  $\mu_{\mathbf{p}}$  and  $\mu_x$ , and variances  $\sigma_{\mathbf{p}}^2$  and  $\sigma_x^2$ . According to the rational expectations assumption, the economic agents know the probability distribution of the economy they live in, i.e., they know the expected values and variances of  $\mathbf{m}$  and  $x$  and the value of  $\gamma$ . But they are unable to observe  $\mathbf{m}$  and  $x$ , they can only observe  $\mathbf{p}$ . Then the best predictor of  $x$  using  $\mathbf{p}$  is the conditional expectation  $E[x|\mathbf{p}]$ .

• a. Assume you are one of these agents and you observe  $\mathbf{p} = p$ . How would you predict  $x$  to be, i.e., what is the value of  $E[x|\mathbf{p} = p]$ ?

ANSWER. It is, according to formula (7.3.18),  $E[x|p = p] = \mu_x + \frac{\text{cov}(x,p)}{\text{var}(p)}(p - E[p])$ . Now  $E[p] = \mu_m + \gamma\mu_x$ ,  $\text{cov}[x, p] = \text{cov}[x, m] + \gamma \text{cov}[x, x] = \gamma\sigma_x^2$ , and  $\text{var}(p) = \sigma_m^2 + \gamma^2\sigma_x^2$ . Therefore

$$(6.7.16) \quad E[x|p = p] = \mu_x + \frac{\gamma\sigma_x^2}{\sigma_m^2 + \gamma^2\sigma_x^2}(p - \mu_m - \gamma\mu_x).$$

□

• b. Define the prediction error  $\varepsilon = x - E[x|p]$ . Compute expected value and variance of  $\varepsilon$ .

ANSWER.

$$(6.7.17) \quad \varepsilon = x - \mu_x - \frac{\gamma\sigma_x^2}{\sigma_m^2 + \gamma^2\sigma_x^2}(p - \mu_m - \gamma\mu_x).$$

This has zero expected value, and its variance is

$$(6.7.18) \quad \text{var}[\varepsilon] = \text{var}[x] + \left(\frac{\gamma\sigma_x^2}{\sigma_m^2 + \gamma^2\sigma_x^2}\right)^2 \text{var}[p] - 2\left(\frac{\gamma\sigma_x^2}{\sigma_m^2 + \gamma^2\sigma_x^2}\right) \text{cov}[x, p] =$$

$$(6.7.19) \quad = \sigma_x^2 + \frac{\gamma^2(\sigma_x^2)^2}{\sigma_m^2 + \gamma^2\sigma_x^2} - 2\frac{\gamma^2(\sigma_x^2)^2}{\sigma_m^2 + \gamma^2\sigma_x^2}$$

$$(6.7.20) \quad = \frac{\sigma_x^2\sigma_m^2}{\sigma_m^2 + \gamma^2\sigma_x^2} = \frac{\sigma_x^2}{1 + \gamma^2\sigma_x^2/\sigma_m^2}.$$

□

• c. In an attempt to fine tune the economy, the central bank increases  $\sigma_m^2$ . Does that increase or decrease  $\text{var}(\varepsilon)$ ?

ANSWER. From (6.7.20) follows that it increases the variance. □

### 6.8. Transformation of Vector Random Variables

In order to obtain the density or probability mass function of a one-to-one transformation of random variables, we have to follow the same 4 steps described in Section 3.6 for a scalar random variable. (1) Determine  $A$ , the range of the new variable, whose density we want to compute; (2) express the old variable, the one whose density/mass function is known, in terms of the new variable, the one whose density or mass function is needed. If that of  $\begin{bmatrix} x \\ y \end{bmatrix}$  is known, set  $\begin{bmatrix} x \\ y \end{bmatrix} = \mathbf{t}(\mathbf{u}, \mathbf{v})$ . Here

$\mathbf{t}$  is a vector-valued function, (i.e., it could be written  $\mathbf{t}(\mathbf{u}, \mathbf{v}) = \begin{bmatrix} \mathbf{q}(\mathbf{u}, \mathbf{v}) \\ \mathbf{r}(\mathbf{u}, \mathbf{v}) \end{bmatrix}$ , but we will use one symbol  $\mathbf{t}$  for this whole transformation), and you have to check that it is one-to-one on  $A$ , i.e.,  $\mathbf{t}(u, v) = \mathbf{t}(u_1, v_1)$  implies  $u = u_1$  and  $v = v_1$  for all  $(u, v)$  and  $(u_1, v_1)$  in  $A$ . (A function for which two different arguments  $(u, v)$  and  $(u_1, v_1)$  give the same function value is called many-to-one.)

If the joint probability distribution of  $\mathbf{x}$  and  $\mathbf{y}$  is described by a probability mass function, then the joint probability mass function of  $\mathbf{u}$  and  $\mathbf{v}$  can simply be obtained

by substituting  $\mathbf{t}$  into the joint probability mass function of  $\mathbf{x}$  and  $\mathbf{y}$  (and it is zero for any values which are not in  $A$ ):

$$(6.8.1) \quad p_{\mathbf{u}, \mathbf{v}}(u, v) = \Pr\left[\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} u \\ v \end{bmatrix}\right] = \Pr[\mathbf{t}(\mathbf{u}, \mathbf{v}) = \mathbf{t}(u, v)] = \Pr\left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{t}(u, v)\right] = p_{\mathbf{x}, \mathbf{y}}(\mathbf{t}(u, v)).$$

The second equal sign is where the condition enters that  $\mathbf{t} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is one-to-one.

If one works with the density function instead of a mass function, one must perform an additional step besides substituting  $\mathbf{t}$ . Since  $\mathbf{t}$  is one-to-one, it follows

$$(6.8.2) \quad \left\{ \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \in dV_{\mathbf{u}, \mathbf{v}} \right\} = \left\{ \mathbf{t}(\mathbf{u}, \mathbf{v}) \in \mathbf{t}(dV)_{\mathbf{x}, \mathbf{y}} \right\}.$$

Therefore

$$(6.8.3) \quad f_{\mathbf{u}, \mathbf{v}}(u, v) |dV_{\mathbf{u}, \mathbf{v}}| = \Pr\left[\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \in dV_{\mathbf{u}, \mathbf{v}}\right] = \Pr[\mathbf{t}(\mathbf{u}, \mathbf{v}) \in \mathbf{t}(dV)_{\mathbf{x}, \mathbf{y}}] = f_{\mathbf{x}, \mathbf{y}}(\mathbf{t}(u, v)) |\mathbf{t}(dV)_{\mathbf{x}, \mathbf{y}}|$$

$$(6.8.4) \quad = f_{\mathbf{x}, \mathbf{y}}(\mathbf{t}(u, v)) \frac{|\mathbf{t}(dV)_{\mathbf{x}, \mathbf{y}}|}{|dV_{\mathbf{u}, \mathbf{v}}|} |dV_{\mathbf{u}, \mathbf{v}}|.$$

The term  $\frac{|\mathbf{t}(dV)_{\mathbf{x}, \mathbf{y}}|}{|dV_{\mathbf{u}, \mathbf{v}}|}$  is the local magnification factor of the transformation. Analytically it is the absolute value  $|J|$  of the Jacobian determinant

$$(6.8.5) \quad J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \begin{vmatrix} \frac{\partial \mathbf{q}}{\partial \mathbf{r}}(u, v) & \frac{\partial \mathbf{q}}{\partial \mathbf{v}}(u, v) \\ \frac{\partial \mathbf{r}}{\partial \mathbf{u}}(u, v) & \frac{\partial \mathbf{r}}{\partial \mathbf{v}}(u, v) \end{vmatrix}.$$

Remember,  $\mathbf{u}, \mathbf{v}$  are the new and  $\mathbf{x}, \mathbf{y}$  the old variables. To compute  $J$  one has to express the old in terms of the new variables. If one expresses the new in terms of the old, one has to take the inverse of the corresponding determinant! The transformation rule for density functions can therefore be summarized as:

$$(\mathbf{x}, \mathbf{y}) = \mathbf{t}(\mathbf{u}, \mathbf{v}) \text{ one-to-one} \Rightarrow f_{\mathbf{u}, \mathbf{v}}(u, v) = f_{\mathbf{x}, \mathbf{y}}(\mathbf{t}(u, v)) |J| \quad \text{where} \quad J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix}$$

PROBLEM 137. Let  $\mathbf{x}$  and  $\mathbf{y}$  be two random variables with joint density function  $f_{\mathbf{x}, \mathbf{y}}(x, y)$ .

• a. 3 points Define  $\mathbf{u} = \mathbf{x} + \mathbf{y}$ . Derive the joint density function of  $\mathbf{u}$  and  $\mathbf{y}$ .

ANSWER. You have to express the “old”  $\mathbf{x}$  and  $\mathbf{y}$  as functions of the “new”  $\mathbf{u}$  and  $\mathbf{y}$ :

$$\begin{matrix} \mathbf{x} = \mathbf{u} - \mathbf{y} \\ \mathbf{y} = \mathbf{y} \end{matrix} \quad \text{or} \quad \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \end{bmatrix} \quad \text{therefore} \quad J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial y} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial y} \end{vmatrix} = \begin{vmatrix} 1 & -1 \\ 0 & 1 \end{vmatrix} = 1.$$

Therefore

$$(6.8.6) \quad f_{\mathbf{u}, \mathbf{y}}(u, y) = f_{\mathbf{x}, \mathbf{y}}(u - y, y).$$

• b. 1 point Derive from this the following formula computing the density function  $f_{\mathbf{u}}(u)$  of the sum  $\mathbf{u} = \mathbf{x} + \mathbf{y}$  from the joint density function  $f_{\mathbf{x},\mathbf{y}}(x, y)$  of  $\mathbf{x}$  and  $\mathbf{y}$ .

$$(6.8.7) \quad f_{\mathbf{u}}(u) = \int_{y=-\infty}^{y=\infty} f_{\mathbf{x},\mathbf{y}}(u - y, y) dy.$$

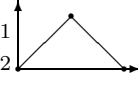
ANSWER. Write down the joint density of  $\mathbf{u}$  and  $\mathbf{y}$  and then integrate  $y$  out, i.e., take its integral over  $y$  from  $-\infty$  to  $+\infty$ :

$$(6.8.8) \quad f_{\mathbf{u}}(u) = \int_{y=-\infty}^{y=\infty} f_{\mathbf{u},\mathbf{y}}(u, y) dy = \int_{y=-\infty}^{y=\infty} f_{\mathbf{x},\mathbf{y}}(u - y, y) dy.$$

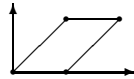
i.e., one integrates over all  $\begin{bmatrix} x \\ y \end{bmatrix}$  with  $x + y = u$ .  $\square$

PROBLEM 138. 6 points Let  $\mathbf{x}$  and  $\mathbf{y}$  be independent and uniformly distributed over the interval  $[0, 1]$ . Compute the density function of  $\mathbf{u} = \mathbf{x} + \mathbf{y}$  and draw its graph. Hint: you may use formula (6.8.7) for the density of the sum of two jointly distributed random variables. An alternative approach would be to first compute the cumulative distribution function  $\Pr[\mathbf{x} + \mathbf{y} \leq u]$  for all  $u$ .

ANSWER. Using equation (6.8.7):

$$(6.8.9) \quad f_{\mathbf{x}+\mathbf{y}}(u) = \int_{-\infty}^{\infty} f_{\mathbf{x},\mathbf{y}}(u - y, y) dy = \begin{cases} u & \text{for } 0 \leq u \leq 1 \\ 2 - u & \text{for } 1 \leq u \leq 2 \\ 0 & \text{otherwise.} \end{cases}$$


To help evaluate this integral, here is the area in  $u, y$ -plane ( $u = x + y$  on the horizontal and  $y$  on the vertical axis) in which  $f_{\mathbf{x},\mathbf{y}}(u - v, v)$  has the value 1:



This is the area between  $(0,0)$ ,  $(1,1)$ ,  $(2,1)$ , and  $(1,0)$ .

One can also show it this way:  $f_{\mathbf{x},\mathbf{y}}(x, y) = 1$  iff  $0 \leq x \leq 1$  and  $0 \leq y \leq 1$ . Now take any fixed  $u$ . It must be between 0 and 2. First assume  $0 \leq u \leq 1$ : then  $f_{\mathbf{x},\mathbf{y}}(u - y, y) = 1$  iff  $0 \leq u - y \leq 1$  and  $0 \leq y \leq 1$  iff  $0 \leq y \leq u$ . Now assume  $1 \leq u \leq 2$ : then  $f_{\mathbf{x},\mathbf{y}}(u - y, y) = 1$  iff  $u - 1 \leq y \leq 1$ .  $\square$

PROBLEM 139. Assume  $\begin{bmatrix} x \\ y \end{bmatrix}$  is uniformly distributed on a round disk around the origin with radius 10.

• a. 4 points Derive the joint density, the marginal density of  $\mathbf{x}$ , and the conditional density of  $\mathbf{y}$  given  $\mathbf{x}=x$ .

• b. 3 points Now let us go over to polar coordinates  $\mathbf{r}$  and  $\phi$ , which satisfy

$$(6.8.10) \quad \begin{aligned} \mathbf{x} &= \mathbf{r} \cos \phi \\ \mathbf{y} &= \mathbf{r} \sin \phi \end{aligned}, \quad \text{i.e., the vector transformation } \mathbf{t} \text{ is } \mathbf{t} \left( \begin{bmatrix} \mathbf{r} \\ \phi \end{bmatrix} \right) = \begin{bmatrix} \mathbf{r} \cos \phi \\ \mathbf{r} \sin \phi \end{bmatrix}.$$

Which region in  $\begin{pmatrix} \mathbf{r} \\ \phi \end{pmatrix}$ -space is necessary to cover  $\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}$ -space? Compute the Jacobian determinant of this transformation. Give an intuitive explanation in terms of local magnification factor of the formula you get. Finally compute the transformed density function.

- c. 1 point Compute  $\text{cov}[\mathbf{x}, \mathbf{y}]$ .
- d. 2 points Compute the conditional variance  $\text{var}[\mathbf{y}|\mathbf{x}=x]$ .
- e. 2 points Are  $\mathbf{x}$  and  $\mathbf{y}$  independent?

PROBLEM 140. [Ame85, pp. 296–7] Assume three transportation choices available: bus, train, and car. If you pick at random a neoclassical individual and ask him or her which utility this person derives from using bus, train, and car the answer will be three numbers  $\mathbf{u}_1(\omega)$ ,  $\mathbf{u}_2(\omega)$ ,  $\mathbf{u}_3(\omega)$ . Here  $\mathbf{u}_1$ ,  $\mathbf{u}_2$ , and  $\mathbf{u}_3$  are assumed to be independent random variables with the following cumulative distribution functions:

$$(6.8.11) \quad \Pr[\mathbf{u}_i \leq u] = F_i(u) = \exp(-\exp(\mu_i - u)), \quad i = 1, 2, 3.$$

I.e., the functional form is the same for all three transportation choices (exp indicates the exponential function); the  $F_i$  only differ by the parameters  $\mu_i$ . The probability distributions are called Type I extreme value distributions, or log Weibull distributions.

Often these kinds of models are set up in such a way that these  $\mu_i$  to depend on the income etc. of the individual, but we assume for this exercise that this distribution applies to the population as a whole.

• a. 1 point Show that the  $F_i$  are indeed cumulative distribution functions, and derive the density functions  $f_i(u)$ .

Individual  $\omega$  likes cars best if and only if his utilities satisfy  $\mathbf{u}_3(\omega) \geq \mathbf{u}_1(\omega)$  and  $\mathbf{u}_3(\omega) \geq \mathbf{u}_2(\omega)$ . Let  $I$  be a function of three arguments such that  $I(\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3)$  is an indicator function of the event that one randomly chooses an individual  $\omega$  who likes cars best, i.e.,

$$(6.8.12) \quad I(\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3) = \begin{cases} 1 & \text{if } \mathbf{u}_1 \leq \mathbf{u}_3 \text{ and } \mathbf{u}_2 \leq \mathbf{u}_3 \\ 0 & \text{otherwise.} \end{cases}$$

Then  $\Pr[\text{car}] = \mathbb{E}[I(\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3)]$ . The following steps have the purpose to compute this probability:

- b. 2 points For any fixed number  $u$ , define  $g(u) = E[I(\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3) | \mathbf{u}_3 = u]$ .

Show that

$$(6.8.13) \quad g(u) = \exp(-\exp(\mu_1 - u) - \exp(\mu_2 - u)).$$

- c. 2 points This here is merely the evaluation of an integral. Show that

$$\begin{aligned} \int_{-\infty}^{+\infty} \exp(-\exp(\mu_1 - u) - \exp(\mu_2 - u) - \exp(\mu_3 - u)) \exp(\mu_3 - u) du &= \\ &= \frac{\exp \mu_3}{\exp \mu_1 + \exp \mu_2 + \exp \mu_3}. \end{aligned}$$

Hint: use substitution rule with  $y = -\exp(\mu_1 - u) - \exp(\mu_2 - u) - \exp(\mu_3 - u)$ .

- d. 1 point Use b and c to show that

$$(6.8.14) \quad \Pr[car] = \frac{\exp \mu_3}{\exp \mu_1 + \exp \mu_2 + \exp \mu_3}.$$

## CHAPTER 7

## The Multivariate Normal Probability Distribution

## 7.1. More About the Univariate Case

By definition,  $z$  is a *standard* normal variable, in symbols,  $z \sim N(0, 1)$ , if it has the density function

$$(7.1.1) \quad f_z(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}.$$

To verify that this is a density function we have to check two conditions. (1) It is everywhere nonnegative. (2) Its integral from  $-\infty$  to  $\infty$  is 1. In order to evaluate this integral, it is easier to work with the independent product of two standard normal variables  $x$  and  $y$ ; their joint density function is  $f_{x,y}(x, y) = \frac{1}{2\pi} e^{-\frac{x^2+y^2}{2}}$ . In order to see that this joint density integrates to 1, go over to polar coordinates  $x = r \cos \phi$ ,  $y = r \sin \phi$ , i.e., compute the joint distribution of  $r$  and  $\phi$  from that of  $x$  and  $y$ : the absolute value of the Jacobian determinant is  $r$ , i.e.,  $dx dy = r dr d\phi$ , therefore

$$(7.1.2) \quad \int_{y=-\infty}^{y=\infty} \int_{x=-\infty}^{x=\infty} \frac{1}{2\pi} e^{-\frac{x^2+y^2}{2}} dx dy = \int_{\phi=0}^{2\pi} \int_{r=0}^{\infty} \frac{1}{2\pi} e^{-\frac{r^2}{2}} r dr d\phi.$$

By substituting  $t = r^2/2$ , therefore  $dt = r dr$ , the inner integral becomes  $-\frac{1}{2\pi} e^{-t} \Big|_0^\infty = \frac{1}{2\pi}$ ; therefore the whole integral is 1. Therefore the product of the integrals of the marginal densities is 1, and since each such marginal integral is positive and they are equal, each of the marginal integrals is 1 too.

**PROBLEM 141.** *6 points* The Gamma function can be defined as  $\Gamma(r) = \int_0^\infty x^{r-1} e^{-x} dx$ . Show that  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ . (Hint: after substituting  $r = 1/2$ , apply the variable transformation  $x = z^2/2$  for nonnegative  $x$  and  $z$  only, and then reduce the resulting integral to the integral over the normal density function.)

**ANSWER.** Then  $dx = z dz$ ,  $\frac{dx}{\sqrt{x}} = dz \sqrt{2}$ . Therefore one can reduce it to the integral over the normal density:

$$(7.1.3) \quad \int_0^\infty \frac{1}{\sqrt{x}} e^{-x} dx = \sqrt{2} \int_0^\infty e^{-z^2/2} dz = \frac{1}{\sqrt{2}} \int_{-\infty}^\infty e^{-z^2/2} dz = \frac{\sqrt{2\pi}}{\sqrt{2}} = \sqrt{\pi}.$$

□

A univariate normal variable with mean  $\mu$  and variance  $\sigma^2$  is a variable  $x$  whose standardized version  $z = \frac{x-\mu}{\sigma} \sim N(0, 1)$ . In this transformation from  $x$  to  $z$ , the Jacobian determinant is  $\frac{dz}{dx} = \frac{1}{\sigma}$ ; therefore the density function of  $x \sim N(\mu, \sigma^2)$  (two notations, the second is perhaps more modern):

$$(7.1.4) \quad f_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = (2\pi\sigma^2)^{-1/2} \exp(-(x-\mu)^2/2\sigma^2).$$

**PROBLEM 142.** *3 points* Given  $n$  independent observations of a Normally distributed variable  $\mathbf{y} \sim N(\mu, 1)$ . Show that the sample mean  $\bar{y}$  is a sufficient statistic for  $\mu$ . Here is a formulation of the factorization theorem for sufficient statistics, which you will need for this question: Given a family of probability density functions  $f_{\mathbf{y}}(y_1, \dots, y_n; \theta)$  defined on  $\mathbb{R}^n$ , which depend on a parameter  $\theta \in \Theta$ . The statistic  $T: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $y_1, \dots, y_n \mapsto T(y_1, \dots, y_n)$  is sufficient for parameter  $\theta$  if and only if there exists a function of two variables  $g: \mathbb{R} \times \Theta \rightarrow \mathbb{R}$ ,  $t, \theta \mapsto g(t; \theta)$ , and a function of  $n$  variables  $h: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $y_1, \dots, y_n \mapsto h(y_1, \dots, y_n)$  so that

$$(7.1.5) \quad f_{\mathbf{y}}(y_1, \dots, y_n; \theta) = g(T(y_1, \dots, y_n); \theta) \cdot h(y_1, \dots, y_n).$$

**ANSWER.** The joint density function can be written (factorization indicated by  $\cdot$ ):

$$(7.1.6) \quad (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - \mu)^2\right) = (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - \bar{y})^2\right) \cdot \exp\left(-\frac{n}{2} (\bar{y} - \mu)^2\right) = h(y_1, \dots, y_n) \cdot g(\bar{y}; \mu).$$

## 7.2. Definition of Multivariate Normal

The multivariate normal distribution is an important family of distributions with very nice properties. But one must be a little careful how to define it. One might naively think a multivariate Normal is a vector random variable each component of which is univariate Normal. But this is not the right definition. Normality of the components is a necessary but not sufficient condition for a multivariate normal vector. If  $\mathbf{u} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$  with both  $\mathbf{x}$  and  $\mathbf{y}$  multivariate normal,  $\mathbf{u}$  is not necessarily multivariate normal.

Here is a recursive definition from which one gets all multivariate normal distributions:

(1) The univariate standard normal  $z$ , considered as a vector with one component, is multivariate normal.

(2) If  $\mathbf{x}$  and  $\mathbf{y}$  are multivariate normal and they are independent, then  $\mathbf{u} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$  is multivariate normal.

(3) If  $\mathbf{y}$  is multivariate normal, and  $\mathbf{A}$  a matrix of constants (which need not be square and is allowed to be singular), and  $\mathbf{b}$  a vector of constants, then  $\mathbf{A}\mathbf{y} + \mathbf{b}$  is multivariate normal. In words: A vector consisting of linear combinations of the same set of multivariate normal variables is again multivariate normal.

For simplicity we will go over now to the bivariate Normal distribution.

### 7.3. Special Case: Bivariate Normal

The following two simple rules allow to obtain all bivariate Normal random variables:

(1) If  $x$  and  $y$  are independent and each of them has a (univariate) normal distribution with mean 0 and the same variance  $\sigma^2$ , then they are bivariate normal. (They would be bivariate normal even if their variances were different and their means not zero, but for the calculations below we will use only this special case, which together with principle (2) is sufficient to get all bivariate normal distributions.)

(2) If  $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$  is bivariate normal and  $\mathbf{P}$  is a  $2 \times 2$  nonrandom matrix and  $\boldsymbol{\mu}$  a nonrandom column vector with two elements, then  $\mathbf{P}\mathbf{x} + \boldsymbol{\mu}$  is bivariate normal as well.

All other properties of bivariate Normal variables can be derived from this.

First let us derive the density function of a bivariate Normal distribution. Write  $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$ .  $x$  and  $y$  are independent  $N(0, \sigma^2)$ . Therefore by principle (1) above the vector  $\mathbf{x}$  is bivariate normal. Take any nonsingular  $2 \times 2$  matrix  $\mathbf{P}$  and a 2 vector  $\boldsymbol{\mu} = \begin{bmatrix} \mu \\ \nu \end{bmatrix}$ , and define  $\begin{bmatrix} u \\ v \end{bmatrix} = \mathbf{u} = \mathbf{P}\mathbf{x} + \boldsymbol{\mu}$ . We need nonsingularity because otherwise the resulting variable would not have a bivariate density; its probability mass would be concentrated on one straight line in the two-dimensional plane. What is the joint density function of  $\mathbf{u}$ ? Since  $\mathbf{P}$  is nonsingular, the transformation is on-to-one, therefore we can apply the transformation theorem for densities. Let us first write down the density function of  $\mathbf{x}$  which we know:

$$(7.3.1) \quad f_{x,y}(x,y) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(x^2 + y^2)\right).$$

For the next step, remember that we have to express the old variable in terms of the new one:  $\mathbf{x} = \mathbf{P}^{-1}(\mathbf{u} - \boldsymbol{\mu})$ . The Jacobian determinant is therefore  $J = \det(\mathbf{P}^{-1})$ . Also notice that, after the substitution  $\begin{bmatrix} x \\ y \end{bmatrix} = \mathbf{P}^{-1} \begin{bmatrix} u - \mu \\ v - \nu \end{bmatrix}$ , the exponent in the joint density function of  $x$  and  $y$  is  $-\frac{1}{2\sigma^2}(x^2 + y^2) = -\frac{1}{2\sigma^2} \begin{bmatrix} x \\ y \end{bmatrix}^\top \begin{bmatrix} x \\ y \end{bmatrix} =$

$-\frac{1}{2\sigma^2} \begin{bmatrix} u - \mu \\ v - \nu \end{bmatrix}^\top \mathbf{P}^{-1\top} \mathbf{P}^{-1} \begin{bmatrix} u - \mu \\ v - \nu \end{bmatrix}$ . Therefore the transformation theorem of densities functions gives

$$(7.3.2) \quad f_{u,v}(u,v) = \frac{1}{2\pi\sigma^2} |\det(\mathbf{P}^{-1})| \exp\left(-\frac{1}{2\sigma^2} \begin{bmatrix} u - \mu \\ v - \nu \end{bmatrix}^\top \mathbf{P}^{-1\top} \mathbf{P}^{-1} \begin{bmatrix} u - \mu \\ v - \nu \end{bmatrix}\right).$$

This expression can be made nicer. Note that the covariance matrix of the transformed variables is  $\mathcal{V}\left[\begin{bmatrix} u \\ v \end{bmatrix}\right] = \sigma^2 \mathbf{P}\mathbf{P}^\top = \sigma^2 \boldsymbol{\Psi}$ , say. Since  $\mathbf{P}^{-1\top} \mathbf{P}^{-1} \mathbf{P}\mathbf{P}^\top = \mathbf{I}$  it follows  $\mathbf{P}^{-1\top} \mathbf{P}^{-1} = \boldsymbol{\Psi}^{-1}$  and  $|\det(\mathbf{P}^{-1})| = 1/\sqrt{|\det(\boldsymbol{\Psi})|}$ , therefore

$$(7.3.3) \quad f_{u,v}(u,v) = \frac{1}{2\pi\sigma^2} \frac{1}{\sqrt{|\det(\boldsymbol{\Psi})|}} \exp\left(-\frac{1}{2\sigma^2} \begin{bmatrix} u - \mu \\ v - \nu \end{bmatrix}^\top \boldsymbol{\Psi}^{-1} \begin{bmatrix} u - \mu \\ v - \nu \end{bmatrix}\right).$$

This is the general formula for the density function of a bivariate normal with nonsingular covariance matrix  $\sigma^2 \boldsymbol{\Psi}$  and mean vector  $\boldsymbol{\mu}$ . One can also use the following notation which is valid for the multivariate Normal variable with  $n$  dimensions, with mean vector  $\boldsymbol{\mu}$  and nonsingular covariance matrix  $\sigma^2 \boldsymbol{\Psi}$ :

$$(7.3.4) \quad f_{\mathbf{x}}(\mathbf{x}) = (2\pi\sigma^2)^{-n/2} (\det \boldsymbol{\Psi})^{-1/2} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Psi}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right).$$

PROBLEM 143. 1 point Show that the matrix product of  $(\mathbf{P}^{-1})^\top \mathbf{P}^{-1}$  and  $\mathbf{P}\mathbf{P}^\top$  is the identity matrix.

PROBLEM 144. 3 points All vectors in this question are  $n \times 1$  column vectors.  $\mathbf{y} = \boldsymbol{\alpha} + \boldsymbol{\varepsilon}$ , where  $\boldsymbol{\alpha}$  is a vector of constants and  $\boldsymbol{\varepsilon}$  is jointly normal with  $\mathcal{E}[\boldsymbol{\varepsilon}] = \mathbf{o}$ . Often, the covariance matrix  $\mathcal{V}[\boldsymbol{\varepsilon}]$  is not given directly, but a  $n \times n$  nonsingular matrix  $\mathbf{T}$  is known which has the property that the covariance matrix of  $\mathbf{T}\boldsymbol{\varepsilon}$  is  $\sigma^2$  times the  $n \times n$  unit matrix, i.e.,

$$(7.3.5) \quad \mathcal{V}[\mathbf{T}\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}_n.$$

Show that in this case the density function of  $\mathbf{y}$  is

$$(7.3.6) \quad f_{\mathbf{y}}(\mathbf{y}) = (2\pi\sigma^2)^{-n/2} |\det(\mathbf{T})| \exp\left(-\frac{1}{2\sigma^2} (\mathbf{T}(\mathbf{y} - \boldsymbol{\alpha}))^\top \mathbf{T}(\mathbf{y} - \boldsymbol{\alpha})\right).$$

Hint: define  $\mathbf{z} = \mathbf{T}\boldsymbol{\varepsilon}$ , write down the density function of  $\mathbf{z}$ , and make a transformation between  $\mathbf{z}$  and  $\mathbf{y}$ .

ANSWER. Since  $\mathcal{E}[\mathbf{z}] = \mathbf{o}$  and  $\mathcal{V}[\mathbf{z}] = \sigma^2 \mathbf{I}_n$ , its density function is  $(2\pi\sigma^2)^{-n/2} \exp(-\mathbf{z}^\top \mathbf{z}/2\sigma^2)$ . Now express  $\mathbf{z}$ , whose density we know, as a function of  $\mathbf{y}$ , whose density function we want to know.



$\mathbf{z} = \mathbf{T}(\mathbf{y} - \boldsymbol{\alpha})$  or

$$(7.3.7) \quad z_1 = t_{11}(y_1 - \alpha_1) + t_{12}(y_2 - \alpha_2) + \cdots + t_{1n}(y_n - \alpha_n)$$

$$(7.3.8) \quad \vdots$$

$$(7.3.9) \quad z_n = t_{n1}(y_1 - \alpha_1) + t_{n2}(y_2 - \alpha_2) + \cdots + t_{nn}(y_n - \alpha_n)$$

therefore the Jacobian determinant is  $\det(\mathbf{T})$ . This gives the result.  $\square$

### 7.3.1. Most Natural Form of Bivariate Normal Density.

PROBLEM 145. *In this exercise we will write the bivariate normal density in its most natural form. For this we set the multiplicative “nuisance parameter”  $\sigma^2 = 1$ , i.e., write the covariance matrix as  $\boldsymbol{\Psi}$  instead of  $\sigma^2\boldsymbol{\Psi}$ .*

• a. 1 point Write the covariance matrix  $\boldsymbol{\Psi} = \mathcal{V}\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}$  in terms of the standard deviations  $\sigma_u$  and  $\sigma_v$  and the correlation coefficient  $\rho$ .

• b. 1 point Show that the inverse of a  $2 \times 2$  matrix has the following form:

$$(7.3.10) \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}.$$

• c. 2 points Show that

$$(7.3.11) \quad \mathbf{q}^2 = [\mathbf{u} - \boldsymbol{\mu} \quad \mathbf{v} - \nu] \boldsymbol{\Psi}^{-1} \begin{bmatrix} \mathbf{u} - \boldsymbol{\mu} \\ \mathbf{v} - \nu \end{bmatrix}$$

$$(7.3.12) \quad = \frac{1}{1 - \rho^2} \left( \frac{(\mathbf{u} - \boldsymbol{\mu})^2}{\sigma_u^2} - 2\rho \frac{\mathbf{u} - \boldsymbol{\mu}}{\sigma_u} \frac{\mathbf{v} - \nu}{\sigma_v} + \frac{(\mathbf{v} - \nu)^2}{\sigma_v^2} \right).$$

• d. 2 points Show the following quadratic decomposition:

$$(7.3.13) \quad \mathbf{q}^2 = \frac{(\mathbf{u} - \boldsymbol{\mu})^2}{\sigma_u^2} + \frac{1}{(1 - \rho^2)\sigma_v^2} \left( \mathbf{v} - \nu - \rho \frac{\sigma_v}{\sigma_u} (\mathbf{u} - \boldsymbol{\mu}) \right)^2.$$

• e. 1 point Show that (7.3.13) can also be written in the form

$$(7.3.14) \quad \mathbf{q}^2 = \frac{(\mathbf{u} - \boldsymbol{\mu})^2}{\sigma_u^2} + \frac{\sigma_u^2}{\sigma_u^2\sigma_v^2 - (\sigma_{uv})^2} \left( \mathbf{v} - \nu - \frac{\sigma_{uv}}{\sigma_u^2} (\mathbf{u} - \boldsymbol{\mu}) \right)^2.$$

• f. 1 point Show that  $d = \sqrt{\det \boldsymbol{\Psi}}$  can be split up, not additively but multiplicatively, as follows:  $d = \sigma_u \cdot \sigma_v \sqrt{1 - \rho^2}$ .

• g. 1 point Using these decompositions of  $d$  and  $\mathbf{q}^2$ , show that the density function  $f_{\mathbf{u},\mathbf{v}}(u, v)$  reads

$$(7.3.15) \quad \frac{1}{\sqrt{2\pi\sigma_u^2}} \exp\left(-\frac{(u - \mu)^2}{2\sigma_u^2}\right) \cdot \frac{1}{\sqrt{2\pi\sigma_v^2}\sqrt{1 - \rho^2}} \exp\left(-\frac{((v - \nu) - \rho \frac{\sigma_v}{\sigma_u}(u - \mu))^2}{2(1 - \rho^2)\sigma_v^2}\right).$$

The second factor in (7.3.15) is the density of a  $\mathbf{N}(\rho \frac{\sigma_v}{\sigma_u} u, (1 - \rho^2)\sigma_v^2)$  evaluated at  $v$ , and the first factor does not depend on  $v$ . Therefore if I integrate  $v$  out I get the marginal density of  $\mathbf{u}$ , this simply gives me the first factor. The conditional density of  $\mathbf{v}$  given  $\mathbf{u} = u$  is the joint divided by the marginal, i.e., it is the second factor. In other words, by completing the square we wrote the joint density function in its natural form as the product of a marginal and a conditional density function:  $f_{\mathbf{u},\mathbf{v}}(u, v) = f_{\mathbf{u}}(u) \cdot f_{\mathbf{v}|\mathbf{u}}(v; u)$ .

From this decomposition one can draw the following conclusions:

- $\mathbf{u} \sim \mathbf{N}(0, \sigma_u^2)$  is normal and, by symmetry,  $\mathbf{v}$  is normal as well. Note that  $\mathbf{u}$  (or  $\mathbf{v}$ ) can be chosen to be any nonzero linear combination of  $\mathbf{x}$  and  $\mathbf{y}$ . A nonzero linear transformation of independent standard normal variables is therefore univariate normal.
- If  $\rho = 0$  then the joint density function is the product of two independent univariate normal density functions. In other words, if the variables are jointly normal, then they are independent whenever they are uncorrelated. In general distributions only the reverse is true.
- The conditional density of  $\mathbf{v}$  conditionally on  $\mathbf{u} = u$  is the second term on the rhs of (7.3.15), i.e., it is normal too.
- The conditional mean is

$$(7.3.16) \quad \mathbf{E}[\mathbf{v}|\mathbf{u} = u] = \rho \frac{\sigma_v}{\sigma_u} u,$$

i.e., it is a linear function of  $u$ . If the (unconditional) means are not zero then the conditional mean is

$$(7.3.17) \quad \mathbf{E}[\mathbf{v}|\mathbf{u} = u] = \mu_v + \rho \frac{\sigma_v}{\sigma_u} (u - \mu_u).$$

Since  $\rho = \frac{\text{cov}[\mathbf{u}, \mathbf{v}]}{\sigma_u \sigma_v}$ , (7.3.17) can also be written as follows:

$$(7.3.18) \quad \mathbf{E}[\mathbf{v}|\mathbf{u} = u] = \mathbf{E}[\mathbf{v}] + \frac{\text{cov}[\mathbf{u}, \mathbf{v}]}{\text{var}[\mathbf{u}]} (u - \mathbf{E}[\mathbf{u}])$$

- The conditional variance is the same whatever value of  $u$  was chosen: value is

$$(7.3.19) \quad \text{var}[\mathbf{v}|\mathbf{u} = u] = \sigma_v^2(1 - \rho^2),$$

which can also be written as

$$(7.3.20) \quad \text{var}[v|u = u] = \text{var}[v] - \frac{(\text{cov}[u, v])^2}{\text{var}[u]}.$$

We did this in such detail because any bivariate normal with zero mean has this form. A multivariate normal distribution is determined by its means and variances and covariances (or correlations coefficients). If the means are not zero, then the densities merely differ from the above by an additive constant in the arguments, i.e., if one needs formulas for nonzero mean, one has to replace  $u$  and  $v$  in the above equations by  $u - \mu_u$  and  $v - \mu_v$ .  $du$  and  $dv$  remain the same, because the Jacobian of the translation  $u \mapsto u - \mu_u$ ,  $v \mapsto v - \mu_v$  is 1. While the univariate normal was determined by mean and standard deviation, the bivariate normal is determined by the two means  $\mu_u$  and  $\mu_v$ , the two standard deviations  $\sigma_u$  and  $\sigma_v$ , and the correlation coefficient  $\rho$ .

### 7.3.2. Level Lines of the Normal Density.

PROBLEM 146. 8 points Define the angle  $\delta = \arccos(\rho)$ , i.e.,  $\rho = \cos \delta$ . In terms of  $\delta$ , the covariance matrix (??) has the form

$$(7.3.21) \quad \Psi = \begin{bmatrix} \sigma_u^2 & \sigma_u \sigma_v \cos \delta \\ \sigma_u \sigma_v \cos \delta & \sigma_v^2 \end{bmatrix}$$

Show that for all  $\phi$ , the vector

$$(7.3.22) \quad \mathbf{x} = \begin{bmatrix} r \sigma_u \cos \phi \\ r \sigma_v \cos(\phi + \delta) \end{bmatrix}$$

satisfies  $\mathbf{x}^\top \Psi^{-1} \mathbf{x} = r^2$ . The opposite holds too, all vectors  $\mathbf{x}$  satisfying  $\mathbf{x}^\top \Psi^{-1} \mathbf{x} = r^2$  can be written in the form (7.3.22) for some  $\phi$ , but I am not asking to prove this. This formula can be used to draw level lines of the bivariate Normal density and confidence ellipses, more details in (??).

PROBLEM 147. The ellipse in Figure 1 contains all the points  $x, y$  for which

$$(7.3.23) \quad [x-1 \quad y-1] \begin{bmatrix} 0.5 & -0.25 \\ -0.25 & 1 \end{bmatrix}^{-1} \begin{bmatrix} x-1 \\ y-1 \end{bmatrix} \leq 6$$

• a. 3 points Compute the probability that a random variable

$$(7.3.24) \quad \begin{bmatrix} x \\ y \end{bmatrix} \sim N\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0.5 & -0.25 \\ -0.25 & 1 \end{bmatrix}\right)$$

falls into this ellipse. Hint: you should apply equation (7.4.9). Then you will have to look up the values of a  $\chi^2$  distribution in a table, or use your statistics software to get it.

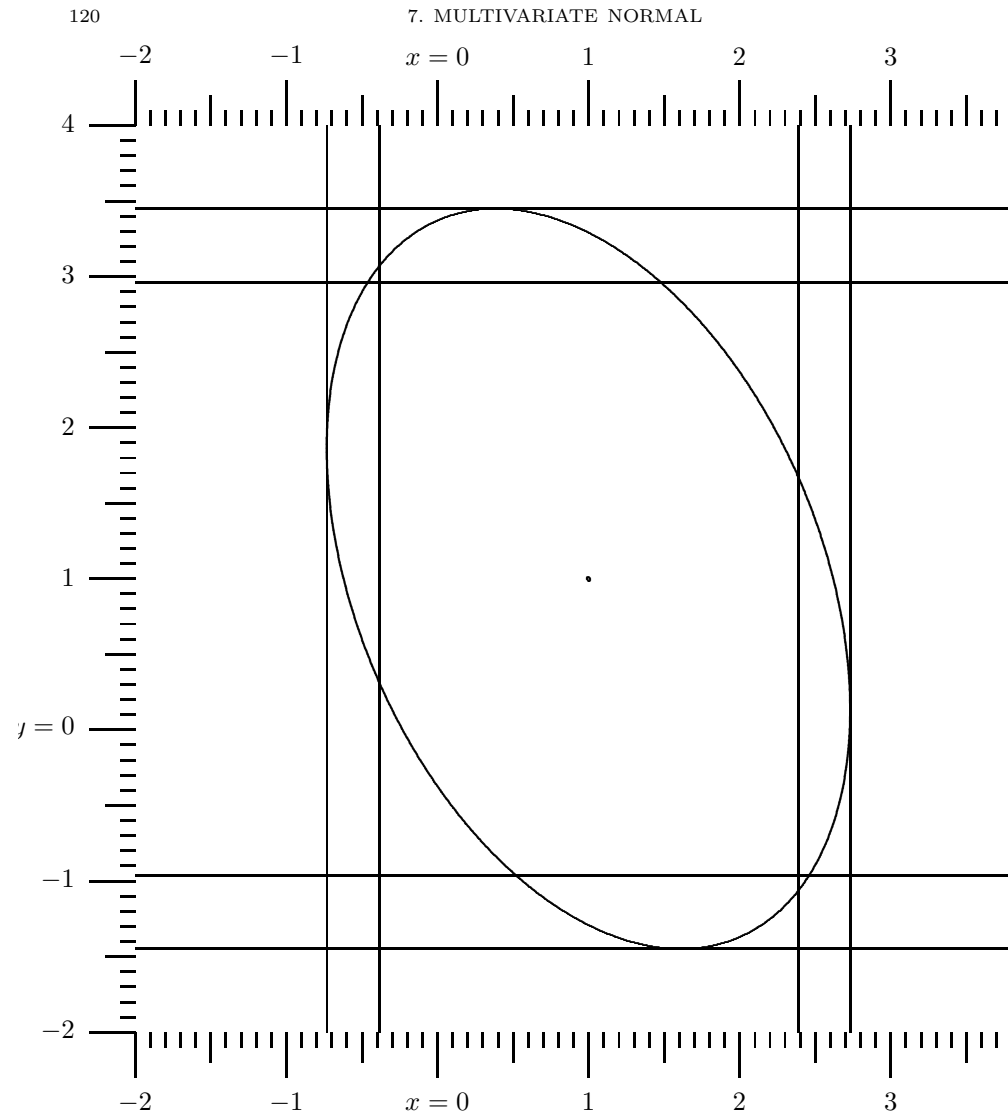


FIGURE 1. Level Line for Normal Density

• b. 1 point Compute the standard deviations of  $x$  and  $y$ , and the correlation coefficient  $\text{corr}(x, y)$

• c. 2 points The vertical tangents to the ellipse in Figure 1 are at the locations  $x = 1 \pm \sqrt{3}$ . What is the probability that  $\begin{bmatrix} x \\ y \end{bmatrix}$  falls between these two vertical tangents?

• d. 1 point The horizontal tangents are at the locations  $y = 1 \pm \sqrt{6}$ . What is the probability that  $\begin{bmatrix} x \\ y \end{bmatrix}$  falls between the horizontal tangents?

• e. 1 point Now take an arbitrary linear combination  $u = ax + by$ . Write down its mean and its standard deviation.

• f. 1 point Show that the set of realizations  $x, y$  for which  $u$  lies less than  $\sqrt{6}$  standard deviation away from its mean is

$$(7.3.25) \quad |a(x-1) + b(y-1)| \leq \sqrt{6} \sqrt{a^2 \text{var}[x] + 2ab \text{cov}[x, y] + b^2 \text{var}[y]}.$$

The set of all these points forms a band limited by two parallel lines. What is the probability that  $\begin{bmatrix} x \\ y \end{bmatrix}$  falls between these two lines?

• g. 1 point It is our purpose to show that this band is again tangent to the ellipse. This is easiest if we use matrix notation. Define

$$(7.3.26) \quad \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix} \quad \boldsymbol{\mu} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \boldsymbol{\Psi} = \begin{bmatrix} 0.5 & -0.25 \\ -0.25 & 1 \end{bmatrix} \quad \mathbf{a} = \begin{bmatrix} a \\ b \end{bmatrix}$$

Equation (7.3.23) in matrix notation says: the ellipse contains all the points for which

$$(7.3.27) \quad (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Psi}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq 6.$$

Show that the band defined by inequality (7.3.25) contains all the points for which

$$(7.3.28) \quad \frac{(\mathbf{a}^\top (\mathbf{x} - \boldsymbol{\mu}))^2}{\mathbf{a}^\top \boldsymbol{\Psi} \mathbf{a}} \leq 6.$$

• h. 2 points Inequality (7.3.28) can also be written as:

$$(7.3.29) \quad (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{a} (\mathbf{a}^\top \boldsymbol{\Psi} \mathbf{a})^{-1} \mathbf{a}^\top (\mathbf{x} - \boldsymbol{\mu}) \leq 6$$

or alternatively

$$(7.3.30) \quad [x-1 \quad y-1] \begin{bmatrix} a \\ b \end{bmatrix} \left( [a \quad b] \boldsymbol{\Psi}^{-1} \begin{bmatrix} a \\ b \end{bmatrix} \right)^{-1} \begin{bmatrix} x-1 \\ y-1 \end{bmatrix} [a \quad b] \leq 6.$$

Show that the matrix

$$(7.3.31) \quad \boldsymbol{\Omega} = \boldsymbol{\Psi}^{-1} - \mathbf{a} (\mathbf{a}^\top \boldsymbol{\Psi} \mathbf{a})^{-1} \mathbf{a}^\top$$

satisfies  $\boldsymbol{\Omega} \boldsymbol{\Psi} \boldsymbol{\Omega} = \boldsymbol{\Omega}$ . Derive from this that  $\boldsymbol{\Omega}$  is nonnegative definite. Hint: you may use, without proof, that any symmetric matrix is nonnegative definite if and only if it can be written in the form  $\mathbf{R} \mathbf{R}^\top$ .

• i. 1 point As an aside: Show that  $\boldsymbol{\Omega} \boldsymbol{\Psi} \mathbf{a} = \mathbf{o}$  and derive from this that  $\boldsymbol{\Omega}$  is not positive definite but only nonnegative definite.

• j. 1 point Show that the following inequality holds for all  $\mathbf{x} - \boldsymbol{\mu}$ ,

$$(7.3.32) \quad (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Psi}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \geq (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{a} (\mathbf{a}^\top \boldsymbol{\Psi} \mathbf{a})^{-1} \mathbf{a}^\top (\mathbf{x} - \boldsymbol{\mu}).$$

In other words, if  $\mathbf{x}$  lies in the ellipse then it also lies in each band. I.e., the ellipse is contained in the intersection of all the bands.

• k. 1 point Show: If  $\mathbf{x} - \boldsymbol{\mu} = \boldsymbol{\Psi} \mathbf{a} \alpha$  with some arbitrary scalar  $\alpha$ , then (7.3.32) is an equality, and if  $\alpha = \pm \sqrt{6 / \mathbf{a}^\top \boldsymbol{\Psi} \mathbf{a}}$ , then both sides in (7.3.32) have the value 6. I.e., the boundary of the ellipse and the boundary lines of the band intersect. Since the ellipse is completely inside the band, this can only be the case if the boundary lines of the band are tangent to the ellipse.

• l. 2 points The vertical lines in Figure 1 which are not tangent to the ellipse delimit a band which, if extended to infinity, has as much probability mass as the ellipse itself. Compute the  $x$ -coordinates of these two lines.

### 7.3.3. Miscellaneous Exercises.

PROBLEM 148. Figure 2 shows the level line for a bivariate Normal density which contains 95% of the probability mass.

• a. 3 points One of the following matrices is the covariance matrix of  $\begin{bmatrix} x \\ y \end{bmatrix}$ .  $\boldsymbol{\Psi}_1$  is  $\begin{bmatrix} 0.62 & -0.56 \\ -0.56 & 1.04 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_2 = \begin{bmatrix} 1.85 & 1.67 \\ 1.67 & 3.12 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_3 = \begin{bmatrix} 0.62 & 0.56 \\ 0.56 & 1.04 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_4 = \begin{bmatrix} 1.85 & -1.67 \\ 1.67 & 3.12 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_5 = \begin{bmatrix} 3.12 & -1.67 \\ -1.67 & 1.85 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_6 = \begin{bmatrix} 1.04 & 0.56 \\ 0.56 & 0.62 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_7 = \begin{bmatrix} 3.12 & 1.67 \\ 1.67 & 1.85 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_8 = \begin{bmatrix} 0.62 & 0.81 \\ 0.81 & 1.04 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_9 = \begin{bmatrix} 3.12 & 1.67 \\ 2.67 & 1.85 \end{bmatrix}$ ,  $\boldsymbol{\Psi}_{10} = \begin{bmatrix} 0.56 & 0.62 \\ 0.62 & -1.04 \end{bmatrix}$ . Which is it? Remember that for a univariate Normal, 95% of the probability mass lie within  $\pm 2$  standard deviations from the mean. If you are not sure, cross out as many of these covariance matrices as possible and write down why you think they should be crossed out.

ANSWER. Covariance matrix must be symmetric, therefore we can cross out 4 and 9. It must also be nonnegative definite (i.e., it must have nonnegative elements in the diagonal), therefore cross out 10, and a nonnegative determinant, therefore cross out 8. Covariance must be positive definite, therefore cross out 1 and 5. Variance in  $x$ -direction is smaller than in  $y$ -direction, therefore cross out 6 and 7. Remains 2 and 3.

Of these it is number 3. By comparison with Figure 1 one can say that the vertical band between 0.4 and 2.6 and the horizontal band between 3 and -1 roughly have the same probability as the ellipse, namely 95%. Since a univariate Normal has 95% of its probability mass in an interval centered around the mean which is 4 standard deviations long, standard deviations must be approximately 0.8 in the horizontal and 1 in the vertical directions.

$\boldsymbol{\Psi}_1$  is negatively correlated;  $\boldsymbol{\Psi}_2$  has the right correlation but is scaled too big;  $\boldsymbol{\Psi}_3$  this is it;  $\boldsymbol{\Psi}_4$  is not symmetric;  $\boldsymbol{\Psi}_5$  negatively correlated, and  $x$  has larger variance than  $y$ ;  $\boldsymbol{\Psi}_6$   $x$  has larger variance

than  $y$ ;  $\Psi_7$  too large,  $x$  has larger variance than  $y$ ;  $\Psi_8$  not positive definite;  $\Psi_9$  not symmetric;  $\Psi_{10}$  not positive definite. □

The next Problem constructs a counterexample which shows that a bivariate distribution, which is not bivariate Normal, can nevertheless have two marginal densities which are univariate Normal.

PROBLEM 149. Let  $x$  and  $y$  be two independent standard normal random variables, and let  $u$  and  $v$  be bivariate normal with mean zero, variances  $\sigma_u^2 = \sigma_v^2 = 1$ , and correlation coefficient  $\rho \neq 0$ . Let  $f_{x,y}$  and  $f_{u,v}$  be the corresponding density functions, i.e.,

$$f_{x,y}(a,b) = \frac{1}{2\pi} \exp\left(-\frac{a^2 + b^2}{2}\right) \quad f_{u,v}(a,b) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-a^2 + b^2 - 2\rho a \frac{b}{2(1-\rho^2)}\right).$$

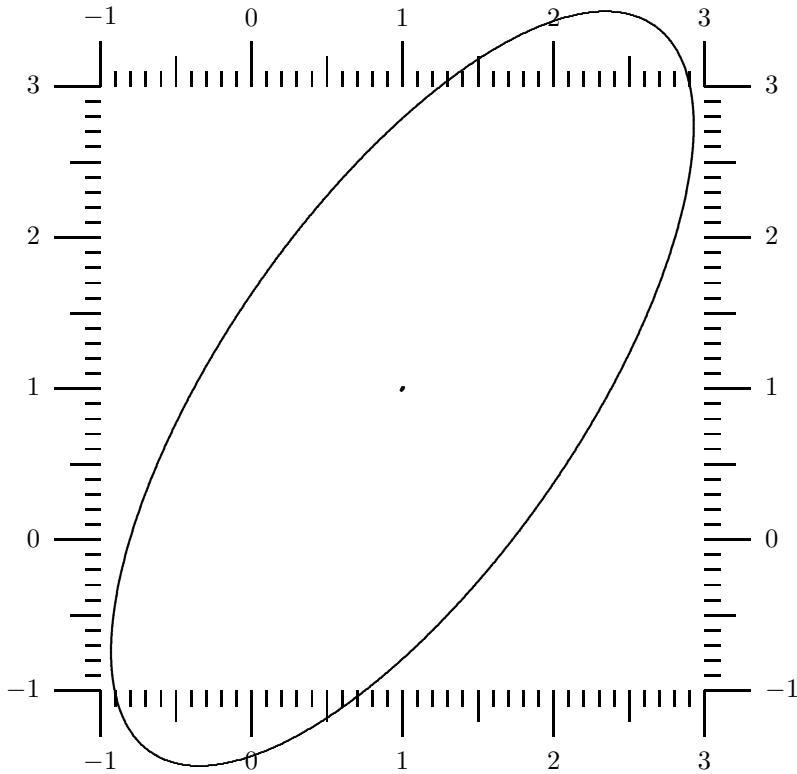


FIGURE 2. Level Line of Bivariate Normal Density, see Problem 148

Assume the random variables  $a$  and  $b$  are defined by the following experiment:  $Y$  flip a fair coin; if it shows head, then you observe  $x$  and  $y$  and give  $a$  the value observed on  $x$ , and  $b$  the value observed of  $y$ . If the coin shows tails, then you observe  $u$  and  $v$  and give  $a$  the value of  $u$ , and  $b$  the value of  $v$ .

- a. Prove that the joint density of  $a$  and  $b$  is

$$(7.3.33) \quad f_{a,b}(a,b) = \frac{1}{2}f_{x,y}(a,b) + \frac{1}{2}f_{u,v}(a,b).$$

Hint: first show the corresponding equation for the cumulative distribution function

ANSWER. Following this hint:

$$(7.3.34) \quad F_{a,b}(a,b) = \Pr[a \leq a \text{ and } b \leq b] =$$

$$(7.3.35) \quad = \Pr[a \leq a \text{ and } b \leq b | \text{head}] \Pr[\text{head}] + \Pr[a \leq a \text{ and } b \leq b | \text{tail}] \Pr[\text{tail}]$$

$$(7.3.36) \quad = F_{x,y}(a,b) \frac{1}{2} + F_{u,v}(a,b) \frac{1}{2}.$$

The density function is the function which, if integrated, gives the above cumulative distribution function.

- b. Show that the marginal distribution of  $a$  and  $b$  each is normal.

ANSWER. You can either argue it out: each of the above marginal distributions is standard normal, but you can also say integrate  $b$  out; for this it is better to use form (7.3.15) for  $f_{u,v}$ , write

$$(7.3.37) \quad f_{u,v}(a,b) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{a^2}{2}\right) \cdot \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp\left(-\frac{(b-\rho a)^2}{2(1-\rho^2)}\right).$$

Then you can see that the marginal is standard normal. Therefore you get a mixture of two distributions each of which is standard normal, therefore it is not really a mixture any more.

- c. Compute the density of  $b$  conditionally on  $a = 0$ . What are its mean and variance? Is it a normal density?

ANSWER.  $F_{b|a}(b; a) = \frac{f_{a,b}(a,b)}{f_a(a)}$ . We don't need it for every  $a$ , only for  $a = 0$ . Since  $f_a(0) = 1/\sqrt{2\pi}$ , therefore

$$(7.3.38) \quad f_{b|a=0}(b) = \sqrt{2\pi} f_{a,b}(0,b) = \frac{1}{2} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{b^2}{2}\right) + \frac{1}{2} \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp\left(-\frac{b^2}{2(1-\rho^2)}\right).$$

It is not normal, it is a mixture of normals with different variances. This has mean zero and variance  $\frac{1}{2}(1 + (1-\rho^2)) = 1 - \frac{1}{2}\rho^2$ .

- d. Are  $a$  and  $b$  jointly normal?

ANSWER. Since the conditional distribution is not normal, they cannot be jointly normal.

PROBLEM 150. This is [HT83, 4.8-6 on p. 263] with variance  $\sigma^2$  instead of 1: Let  $\mathbf{x}$  and  $\mathbf{y}$  be independent normal with mean 0 and variance  $\sigma^2$ . Go over to polar coordinates  $r$  and  $\phi$ , which satisfy

$$(7.3.39) \quad \begin{aligned} \mathbf{x} &= r \cos \phi \\ \mathbf{y} &= r \sin \phi. \end{aligned}$$

- a. 1 point Compute the Jacobian determinant.

ANSWER. Express the variables whose density you know in terms of those whose density you want to know. The Jacobian determinant is

$$(7.3.40) \quad J = \begin{vmatrix} \frac{\partial \mathbf{x}}{\partial r} & \frac{\partial \mathbf{x}}{\partial \phi} \\ \frac{\partial \mathbf{y}}{\partial r} & \frac{\partial \mathbf{y}}{\partial \phi} \end{vmatrix} = \begin{vmatrix} \cos \phi & -r \sin \phi \\ \sin \phi & r \cos \phi \end{vmatrix} = ((\cos \phi)^2 + (\sin \phi)^2)r = r.$$

□

- b. 2 points Find the joint probability density function of  $r$  and  $\phi$ . Also indicate the area in  $(r, \phi)$  space in which it is nonzero.

ANSWER.  $f_{\mathbf{x}, \mathbf{y}}(x, y) = \frac{1}{2\pi\sigma^2} e^{-(x^2+y^2)/2\sigma^2}$ ; therefore  $f_{r, \phi}(r, \phi) = \frac{1}{2\pi\sigma^2} r e^{-r^2/2\sigma^2}$  for  $0 \leq r < \infty$  and  $0 \leq \phi < 2\pi$ . □

- c. 3 points Find the marginal distributions of  $r$  and  $\phi$ . Hint: for one of the integrals it is convenient to make the substitution  $q = r^2/2\sigma^2$ .

ANSWER.  $f_r(r) = \frac{1}{\sigma^2} r e^{-r^2/2\sigma^2}$  for  $0 \leq r < \infty$ , and  $f_\phi(\phi) = \frac{1}{2\pi}$  for  $0 \leq \phi < 2\pi$ . For the latter we need  $\frac{1}{2\pi\sigma^2} \int_0^\infty r e^{-r^2/2\sigma^2} dr = \frac{1}{2\pi}$ , set  $q = r^2/2\sigma^2$ , then  $dq = \frac{1}{\sigma^2} r dr$ , and the integral becomes  $\frac{1}{2\pi} \int_0^\infty e^{-q} dq$ . □

- d. 1 point Are  $r$  and  $\phi$  independent?

ANSWER. Yes, because joint density function is the product of the marginals. □

## 7.4. Multivariate Standard Normal in Higher Dimensions

Here is an important fact about the multivariate normal, which one cannot see in two dimensions: if the partitioned vector  $\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$  is jointly normal, and every component of  $\mathbf{x}$  is independent of every component of  $\mathbf{y}$ , then the vectors  $\mathbf{x}$  and  $\mathbf{y}$  are already independent. Not surprised? You should be, see Problem 125.

Let's go back to the construction scheme at the beginning of this chapter. First we will introduce the multivariate *standard* normal, which one obtains by applying only operations (1) and (2), i.e., it is a vector composed of independent univariate standard normals, and give some properties of it. Then we will go over to the multivariate normal with arbitrary covariance matrix, which is simply an arbitrary linear transformation of the multivariate standard normal. We will always carry the “nuisance parameter”  $\sigma^2$  along.

DEFINITION 7.4.1. The random vector  $\mathbf{z}$  is said to have a multivariate standard normal distribution with variance  $\sigma^2$ , written as  $\mathbf{z} \sim \mathbf{N}(\mathbf{o}, \sigma^2 \mathbf{I})$ , if each element  $z_i$  is a standard normal with same variance  $\sigma^2$ , and all elements are mutually independent of each other. (Note that this definition of the standard normal is a little broader than the usual one; the usual one requires that  $\sigma^2 = 1$ .)

The density function of a multivariate standard normal  $\mathbf{z}$  is therefore the product of the univariate densities, which gives  $f_{\mathbf{z}}(\mathbf{z}) = (2\pi\sigma^2)^{-n/2} \exp(-\mathbf{z}^\top \mathbf{z}/2\sigma^2)$ .

The following property of the multivariate standard normal distributions is basic.

THEOREM 7.4.2. Let  $\mathbf{z}$  be multivariate standard normal  $p$ -vector with variance  $\sigma^2$ , and let  $\mathbf{P}$  be a  $m \times p$  matrix with  $\mathbf{P}\mathbf{P}^\top = \mathbf{I}$ . Then  $\mathbf{x} = \mathbf{P}\mathbf{z}$  is a multivariate standard normal  $m$ -vector with the same variance  $\sigma^2$ , and  $\mathbf{z}^\top \mathbf{z} - \mathbf{x}^\top \mathbf{x} \sim \sigma^2 \chi_{p-m}^2$  independent of  $\mathbf{x}$ .

PROOF.  $\mathbf{P}\mathbf{P}^\top = \mathbf{I}$  means all rows are orthonormal. If  $\mathbf{P}$  is not square, must therefore have more columns than rows, and one can add more rows to get an orthogonal square matrix, call it  $\mathbf{T} = \begin{bmatrix} \mathbf{P} \\ \mathbf{Q} \end{bmatrix}$ . Define  $\mathbf{y} = \mathbf{T}\mathbf{z}$ , i.e.,  $\mathbf{z} = \mathbf{T}^\top \mathbf{y}$ . Then

$\mathbf{z}^\top \mathbf{z} = \mathbf{y}^\top \mathbf{T}\mathbf{T}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{y}$ , and the Jacobian of the transformation from  $\mathbf{y}$  to  $\mathbf{z}$  has absolute value one. Therefore the density function of  $\mathbf{y}$  is  $(2\pi\sigma^2)^{-n/2} \exp(-\mathbf{y}^\top \mathbf{y}/2\sigma^2)$  which means  $\mathbf{y}$  is standard normal as well. In other words, every  $\mathbf{y}_i$  is univariate standard normal with same variance  $\sigma^2$  and  $\mathbf{y}_i$  is independent of  $\mathbf{y}_j$  for  $i \neq j$ . Therefore also any subvector of  $\mathbf{y}$ , such as  $\mathbf{x}$ , is standard normal. Since  $\mathbf{z}^\top \mathbf{z} - \mathbf{x}^\top \mathbf{x} = \mathbf{y}^\top \mathbf{y} - \mathbf{x}^\top \mathbf{x}$  is the sum of the squares of those elements of  $\mathbf{y}$  which are not in  $\mathbf{x}$ , it follows that  $\mathbf{z}^\top \mathbf{z} - \mathbf{x}^\top \mathbf{x}$  is an independent  $\sigma^2 \chi_{p-m}^2$ .

PROBLEM 151. Show that the moment generating function of a multivariate standard normal with variance  $\sigma^2$  is  $m_{\mathbf{z}}(\mathbf{t}) = \mathcal{E}[\exp(\mathbf{t}^\top \mathbf{z})] = \exp(\sigma^2 \mathbf{t}^\top \mathbf{t}/2)$ .

ANSWER. Proof: The moment generating function is defined as

$$(7.4.1) \quad m_{\mathbf{z}}(\mathbf{t}) = \mathcal{E}[\exp(\mathbf{t}^\top \mathbf{z})]$$

$$(7.4.2) \quad = (2\pi\sigma^2)^{n/2} \int \cdots \int \exp(-\frac{1}{2\sigma^2} \mathbf{z}^\top \mathbf{z}) \exp(\mathbf{t}^\top \mathbf{z}) dz_1 \cdots dz_n$$

$$(7.4.3) \quad = (2\pi\sigma^2)^{n/2} \int \cdots \int \exp(-\frac{1}{2\sigma^2} (\mathbf{z} - \sigma^2 \mathbf{t})^\top (\mathbf{z} - \sigma^2 \mathbf{t}) + \frac{\sigma^2}{2} \mathbf{t}^\top \mathbf{t}) dz_1 \cdots dz_n$$

$$(7.4.4) \quad = \exp(\frac{\sigma^2}{2} \mathbf{t}^\top \mathbf{t}) \quad \text{since first part of integrand is density function.}$$

THEOREM 7.4.3. Let  $\mathbf{z} \sim \mathbf{N}(\mathbf{o}, \sigma^2 \mathbf{I})$ , and  $\mathbf{P}$  symmetric and of rank  $r$ . A necessary and sufficient condition for  $\mathbf{q} = \mathbf{z}^\top \mathbf{P}\mathbf{z}$  to have a  $\sigma^2 \chi^2$  distribution is  $\mathbf{P}^2 = \mathbf{0}$ . In this case, the  $\chi^2$  has  $r$  degrees of freedom.

Proof of sufficiency: If  $\mathbf{P}^2 = \mathbf{P}$  with rank  $r$ , then a matrix  $\mathbf{T}$  exists with  $\mathbf{P} = \mathbf{T}^\top \mathbf{T}$  and  $\mathbf{T}\mathbf{T}^\top = \mathbf{I}$ . Define  $\mathbf{x} = \mathbf{T}\mathbf{z}$ ; it is standard normal by theorem 7.4.2. Therefore  $\mathbf{q} = \mathbf{z}^\top \mathbf{T}^\top \mathbf{T}\mathbf{z} = \sum_{i=1}^r x_i^2$ .

Proof of necessity by construction of the moment generating function of  $\mathbf{q} = \mathbf{z}^\top \mathbf{P}\mathbf{z}$  for arbitrary symmetric  $\mathbf{P}$  with rank  $r$ . Since  $\mathbf{P}$  is symmetric, there exists a  $\mathbf{T}$  with  $\mathbf{T}\mathbf{T}^\top = \mathbf{I}_r$  and  $\mathbf{P} = \mathbf{T}^\top \mathbf{\Lambda}\mathbf{T}$  where  $\mathbf{\Lambda}$  is a nonsingular diagonal matrix, write it  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_r)$ . Therefore  $\mathbf{q} = \mathbf{z}^\top \mathbf{T}^\top \mathbf{\Lambda}\mathbf{T}\mathbf{z} = \mathbf{x}^\top \mathbf{\Lambda}\mathbf{x} = \sum_{i=1}^r \lambda_i x_i^2$  where  $\mathbf{x} = \mathbf{T}\mathbf{z} \sim \mathbf{N}(\mathbf{o}, \sigma^2 \mathbf{I}_r)$ . Therefore the moment generating function

$$(7.4.5) \quad \mathbb{E}[\exp(\mathbf{q}t)] = \mathbb{E}[\exp(t \sum_{i=1}^r \lambda_i x_i^2)]$$

$$(7.4.6) \quad = \mathbb{E}[\exp(t\lambda_1 x_1^2)] \cdots \mathbb{E}[\exp(t\lambda_r x_r^2)]$$

$$(7.4.7) \quad = (1 - 2\lambda_1 \sigma^2 t)^{-1/2} \cdots (1 - 2\lambda_r \sigma^2 t)^{-1/2}.$$

By assumption this is equal to  $(1 - 2\sigma^2 t)^{-k/2}$  with some integer  $k \geq 1$ . Taking squares and inverses one obtains

$$(7.4.8) \quad (1 - 2\lambda_1 \sigma^2 t) \cdots (1 - 2\lambda_r \sigma^2 t) = (1 - 2\sigma^2 t)^k.$$

Since the  $\lambda_i \neq 0$ , one obtains  $\lambda_i = 1$  by uniqueness of the polynomial roots. Furthermore, this also implies  $r = k$ .

From Theorem 7.4.3 one can derive a characterization of all the quadratic forms of multivariate normal variables with arbitrary covariance matrices that are  $\chi^2$ 's. Assume  $\mathbf{y}$  is a multivariate normal vector random variable with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\sigma^2 \boldsymbol{\Psi}$ , and  $\boldsymbol{\Omega}$  is a symmetric nonnegative definite matrix. Then  $(\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Omega}(\mathbf{y} - \boldsymbol{\mu}) \sim \sigma^2 \chi_k^2$  iff

$$(7.4.9) \quad \boldsymbol{\Psi}\boldsymbol{\Omega}\boldsymbol{\Psi}\boldsymbol{\Omega}\boldsymbol{\Psi} = \boldsymbol{\Psi}\boldsymbol{\Omega}\boldsymbol{\Psi},$$

and  $k$  is the rank of  $\boldsymbol{\Psi}\boldsymbol{\Omega}$ .

Here are the three best known special cases (with examples):

- $\boldsymbol{\Psi} = \mathbf{I}$  (the identity matrix) and  $\boldsymbol{\Omega}^2 = \boldsymbol{\Omega}$ , i.e., the case of theorem 7.4.3. This is the reason why the minimum value of the SSE has a  $\sigma^2 \chi^2$  distribution, see (27.0.10).
- $\boldsymbol{\Psi}$  nonsingular and  $\boldsymbol{\Omega} = \boldsymbol{\Psi}^{-1}$ . The quadratic form in the exponent of the normal density function is therefore a  $\chi^2$ ; one needs therefore the  $\chi^2$  to compute the probability that the realization of a Normal is in a given equidensity-ellipse (Problem 147).
- $\boldsymbol{\Psi}$  singular and  $\boldsymbol{\Omega} = \boldsymbol{\Psi}^-$ , its g-inverse. The multinomial distribution has a singular covariance matrix, and equation (??) gives a convenient g-inverse which enters the equation for Pearson's goodness of fit test.

Here are, without proof, two more useful theorems about the standard normal:

**THEOREM 7.4.4.** *Let  $\mathbf{x}$  a multivariate standard normal. Then  $\mathbf{x}^\top \mathbf{P}\mathbf{x}$  is independent of  $\mathbf{x}^\top \mathbf{Q}\mathbf{x}$  if and only if  $\mathbf{P}\mathbf{Q} = \mathbf{O}$ .*

This is called Craig's theorem, although Craig's proof in [Cra43] is incorrect. Kshirsagar [Ksh19, p. 41] describes the correct proof; he and Seber [Seb77] give Lancaster's book [Lan69] as basic reference. Seber [Seb77] gives a proof which is only valid if the two quadratic forms are  $\chi^2$ .

The next theorem is known as James's theorem, it is a stronger version of Cochran's theorem. It is from Kshirsagar [Ksh19, p. 41].

**THEOREM 7.4.5.** *Let  $\mathbf{x}$  be  $p$ -variate standard normal with variance  $\sigma^2$ , and  $\mathbf{x}^\top \mathbf{x} = \sum_{i=1}^k \mathbf{x}^\top \mathbf{P}_i \mathbf{x}$ . Then for the quadratic forms  $\mathbf{x}^\top \mathbf{P}_i \mathbf{x}$  to be independently distributed  $\sigma^2 \chi^2$ , any one of the following three equivalent conditions is necessary and sufficient*

$$(7.4.10) \quad \mathbf{P}_i^2 = \mathbf{P}_i \quad \text{for all } i$$

$$(7.4.11) \quad \mathbf{P}_i \mathbf{P}_j = \mathbf{O} \quad i \neq j$$

$$(7.4.12) \quad \sum_{i=1}^k \text{rank}(\mathbf{P}_i) = p$$

## CHAPTER 8

## The Regression Fallacy

Only for the sake of this exercise we will assume that “intelligence” is an innate property of individuals and can be represented by a real number  $z$ . If one picks at random a student entering the U of U, the intelligence of this student is a random variable which we assume to be normally distributed with mean  $\mu$  and standard deviation  $\sigma$ . Also assume every student has to take two intelligence tests, the first at the beginning of his or her studies, the other half a year later. The outcomes of these tests are  $x$  and  $y$ .  $x$  and  $y$  measure the intelligence  $z$  (which is assumed to be the same in both tests) plus a random error  $\varepsilon$  and  $\delta$ , i.e.,

$$(8.0.13) \quad x = z + \varepsilon$$

$$(8.0.14) \quad y = z + \delta$$

Here  $z \sim N(\mu, \tau^2)$ ,  $\varepsilon \sim N(0, \sigma^2)$ , and  $\delta \sim N(0, \sigma^2)$  (i.e., we assume that both errors have the same variance). The three variables  $\varepsilon$ ,  $\delta$ , and  $z$  are independent of each other. Therefore  $x$  and  $y$  are jointly normal.  $\text{var}[x] = \tau^2 + \sigma^2$ ,  $\text{var}[y] = \tau^2 + \sigma^2$ ,  $\text{cov}[x, y] = \text{cov}[z + \varepsilon, z + \delta] = \tau^2 + 0 + 0 + 0 = \tau^2$ . Therefore  $\rho = \frac{\tau^2}{\tau^2 + \sigma^2}$ . The contour lines of the joint density are ellipses with center  $(\mu, \mu)$  whose main axes are the lines  $y = x$  and  $y = -x$  in the  $x, y$ -plane.

Now what is the conditional mean? Since  $\text{var}[x] = \text{var}[y]$ , (7.3.17) gives the line  $E[y|x=x] = \mu + \rho(x - \mu)$ , i.e., it is a line which goes through the center of the ellipses but which is flatter than the line  $x = y$  representing the real underlying linear relationship if there are no errors. Geometrically one can get it as the line which intersects each ellipse exactly where the ellipse is vertical.

Therefore, the parameters of the best prediction of  $y$  on the basis of  $x$  are *not* the parameters of the underlying relationship. Why not? Because not only  $y$  but also  $x$  is subject to errors. Assume you pick an individual by random, and it turns out that his or her first test result is very much higher than the average. Then it is more likely that this is an individual which was lucky in the first exam, and his or her true IQ is lower than the one measured, than that the individual is an Einstein who had a bad day. This is simply because  $z$  is normally distributed, i.e., among the students entering a given University, there are more individuals with lower IQ's than

Einsteins. In order to make a good prediction of the result of the second test one must make allowance for the fact that the individual's IQ is most likely lower than his first score indicated, therefore one will predict the second score to be lower than the first score. The converse is true for individuals who scored lower than average, i.e., in your prediction you will do as if a “regression towards the mean” had taken place.

The next important point to note here is: the “true regression line,” i.e., the prediction line, is uniquely determined by the joint distribution of  $x$  and  $y$ . However, the line representing the underlying relationship can only be determined if one has information in addition to the joint density, i.e., in addition to the observations. E.g., assume the two tests have different standard deviations, which may be the case simply because the second test has more questions and is therefore more accurate. Then the underlying 45° line is no longer one of the main axes of the ellipse! To be more precise, the underlying line can only be identified if one knows the ratio of the variances, or if one knows one of the two variances. Without any knowledge of the variances, the only thing one can say about the underlying line is that it lies between the line predicting  $y$  on the basis of  $x$  and the line predicting  $x$  on the basis of  $y$ .

The name “regression” stems from a confusion between the prediction line and the real underlying relationship. Francis Galton, the cousin of the famous Darwin, measured the height of fathers and sons, and concluded from his evidence that the heights of sons tended to be closer to the average height than the height of the fathers, a purported law of “regression towards the mean.” Problem 152 illustrates this:

**PROBLEM 152.** *The evaluation of two intelligence tests, one at the beginning of the semester, one at the end, gives the following disturbing outcome: While the underlying intelligence during the first test was  $z \sim N(100, 20)$ , it changed between the first and second test due to the learning experience at the university. If  $w$  is the intelligence of each student at the second test, it is connected to his intelligence at the first test by the formula  $w = 0.5z + 50$ , i.e., those students with intelligence below 100 gained, but those students with intelligence above 100 lost. (The errors of both intelligence tests are normally distributed with expected value zero, and the variance of the first intelligence test was 5, and that of the second test, which had more questions, was 4. As usual, the errors are independent of each other and of the actual intelligence.)*

• a. 3 points *If  $x$  and  $y$  are the outcomes of the first and second intelligence test, compute  $E[x]$ ,  $E[y]$ ,  $\text{var}[x]$ ,  $\text{var}[y]$ , and the correlation coefficient  $\rho = \text{corr}[x, y]$ . Figure 1 shows an equi-density line of their joint distribution; 95% of the probability mass of the test results are inside this ellipse. Draw the line  $w = 0.5z + 50$  in Figure 1.*

ANSWER. We know  $z \sim N(100, 20)$ ;  $w = 0.5z + 50$ ;  $x = z + \varepsilon$ ;  $\varepsilon \sim N(0, 4)$ ;  $y = w + \delta$ ;  $\delta \sim N(0, 5)$ ; therefore  $E[x] = 100$ ;  $E[y] = 100$ ;  $\text{var}[x] = 20 + 5 = 25$ ;  $\text{var}[y] = 5 + 4 = 9$ ;  $\text{cov}[x, y] = 10$ ;  $\text{corr}[x, y] = 10/15 = 2/3$ . In matrix notation

$$(8.0.15) \quad \begin{bmatrix} x \\ y \end{bmatrix} \sim N \left[ \begin{bmatrix} 100 \\ 100 \end{bmatrix}, \begin{bmatrix} 25 & 10 \\ 10 & 9 \end{bmatrix} \right]$$

The line  $y = 50 + 0.5x$  goes through the points (80, 90) and (120, 110). □

• b. 4 points Compute  $E[y|x=x]$  and  $E[x|y=y]$ . The first is a linear function of  $x$  and the second a linear function of  $y$ . Draw the two lines representing these linear functions into Figure 1. Use (7.3.18) for this.

ANSWER.

$$(8.0.16) \quad E[y|x=x] = 100 + \frac{10}{25}(x - 100) = 60 + \frac{2}{5}x$$

$$(8.0.17) \quad E[x|y=y] = 100 + \frac{10}{9}(y - 100) = -\frac{100}{9} + \frac{10}{9}y.$$

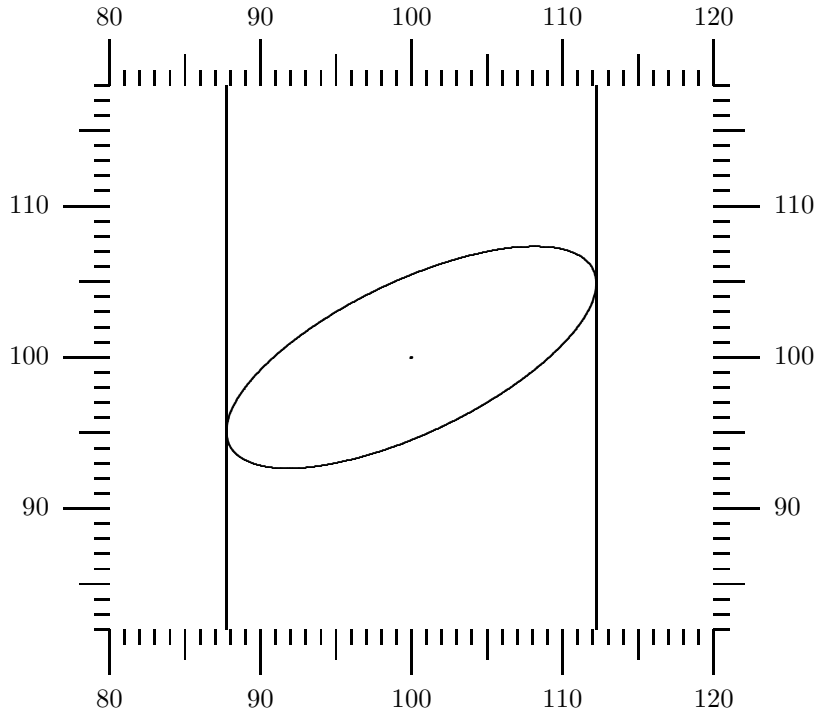


FIGURE 1. Ellipse containing 95% of the probability mass of test results  $x$  and  $y$

The line  $y = E[y|x=x]$  goes through the points (80, 92) and (120, 108) at the edge of Figure 1 and intersects the ellipse where it is vertical. The line  $x = E[x|y=y]$  goes through the points (80, 82) and (120, 118), which are the corner points of Figure 1; it intersects the ellipse where it is horizontal. The two lines intersect in the center of the ellipse, i.e., at the point (100, 100).

• c. 2 points Another researcher says that  $w = \frac{6}{10}z + 40$ ,  $z \sim N(100, \frac{10}{6})$ ,  $\varepsilon \sim N(0, \frac{50}{6})$ ,  $\delta \sim N(0, 3)$ . Is this compatible with the data?

ANSWER. Yes, it is compatible:  $E[x] = E[z] + E[\varepsilon] = 100$ ;  $E[y] = E[w] + E[\delta] = \frac{6}{10}100 + 40 = 100$ ;  $\text{var}[x] = \frac{100}{6} + \frac{50}{6} = 25$ ;  $\text{var}[y] = (\frac{6}{10})^2 \text{var}[z] + \text{var}[\delta] = \frac{63}{100} \frac{100}{6} + 3 = 9$ ;  $\text{cov}[x, y] = \frac{6}{10} \text{var}[z] = 10$ .

• d. 4 points A third researcher asserts that the IQ of the students really did not change. He says  $w = z$ ,  $z \sim N(100, 5)$ ,  $\varepsilon \sim N(0, 20)$ ,  $\delta \sim N(0, 4)$ . Is this compatible with the data? Is there unambiguous evidence in the data that the IQ declined?

ANSWER. This is not compatible. This scenario gets everything right except the covariance:  $E[x] = E[z] + E[\varepsilon] = 100$ ;  $E[y] = E[z] + E[\delta] = 100$ ;  $\text{var}[x] = 5 + 20 = 25$ ;  $\text{var}[y] = 5 + 4 = 9$ ;  $\text{cov}[x, y] = 5$ . A scenario in which both tests have same underlying intelligence cannot be found. Since the two conditional expectations are on the same side of the diagonal, the hypothesis that the intelligence did not change between the two tests is not consistent with the joint distribution of  $x$  and  $y$ . The diagonal goes through the points (82, 82) and (118, 118), i.e., it intersects the horizontal boundaries of Figure 1.

We just showed that the parameters of the true underlying relationship cannot be inferred from the data alone if there are errors in both variables. We also showed that this lack of identification is not complete, because one can specify an interval in which in the plim contains the true parameter value.

Chapter ?? has a much more detailed discussion of all this. There we will show that this lack of identification can be removed if more information is available, i.e., if one knows that the two error variances are equal, or if one knows that the regression line has zero intercept, etc. Question 153 shows that in this latter case, the OLS estimates are not consistent, but other estimates exist that are consistent.

PROBLEM 153. [Fri57, chapter 3] According to Friedman's permanent income hypothesis, drawing at random families in a given country and asking them about their income  $y$  and consumption  $c$  can be modeled as the independent observations of two random variables which satisfy

$$(8.0.18) \quad y = y^p + y^t,$$

$$(8.0.19) \quad c = c^p + c^t,$$

$$(8.0.20) \quad c^p = \beta y^p.$$

Here  $y^p$  and  $c^p$  are the permanent and  $y^t$  and  $c^t$  the transitory components of income and consumption. These components are not observed separately, only the



sums  $\mathbf{y}$  and  $\mathbf{c}$  are observed. We assume that the permanent income  $\mathbf{y}^p$  is random, with  $E[\mathbf{y}^p] = \mu \neq 0$  and  $\text{var}[\mathbf{y}^p] = \tau_y^2$ . The transitory components  $\mathbf{y}^t$  and  $\mathbf{c}^t$  are assumed to be independent of each other and of  $\mathbf{y}^p$ , and  $E[\mathbf{y}^t] = 0$ ,  $\text{var}[\mathbf{y}^t] = \sigma_y^2$ ,  $E[\mathbf{c}^t] = 0$ , and  $\text{var}[\mathbf{c}^t] = \sigma_c^2$ . Finally, it is assumed that all variables are normally distributed.

• a. 2 points Given the above information, write down the vector of expected values  $\mathcal{E}[\begin{bmatrix} \mathbf{y} \\ \mathbf{c} \end{bmatrix}]$  and the covariance matrix  $\mathcal{V}[\begin{bmatrix} \mathbf{y} \\ \mathbf{c} \end{bmatrix}]$  in terms of the five unknown parameters of the model  $\mu$ ,  $\beta$ ,  $\tau_y^2$ ,  $\sigma_y^2$ , and  $\sigma_c^2$ .

ANSWER.

$$(8.0.21) \quad \mathcal{E}\left[\begin{bmatrix} \mathbf{y} \\ \mathbf{c} \end{bmatrix}\right] = \begin{bmatrix} \mu \\ \beta\mu \end{bmatrix} \quad \text{and} \quad \mathcal{V}\left[\begin{bmatrix} \mathbf{y} \\ \mathbf{c} \end{bmatrix}\right] = \begin{bmatrix} \tau_y^2 + \sigma_y^2 & \beta\tau_y^2 \\ \beta\tau_y^2 & \beta^2\tau_y^2 + \sigma_c^2 \end{bmatrix}.$$

□

• b. 3 points Assume that you know the true parameter values and you observe a family's actual income  $\mathbf{y}$ . Show that your best guess (minimum mean squared error) of this family's permanent income  $\mathbf{y}^p$  is

$$(8.0.22) \quad \mathbf{y}^{p*} = \frac{\sigma_y^2}{\tau_y^2 + \sigma_y^2} \mu + \frac{\tau_y^2}{\tau_y^2 + \sigma_y^2} \mathbf{y}.$$

Note: here we are guessing income, not yet consumption! Use (7.3.17) for this!

ANSWER. This answer also does the math for part c. The best guess is the conditional mean

$$\begin{aligned} E[\mathbf{y}^p | \mathbf{y} = 22,000] &= E[\mathbf{y}^p] + \frac{\text{cov}[\mathbf{y}^p, \mathbf{y}]}{\text{var}[\mathbf{y}]} (22,000 - E[\mathbf{y}]) \\ &= 12,000 + \frac{16,000,000}{20,000,000} (22,000 - 12,000) = 20,000 \end{aligned}$$

or equivalently

$$\begin{aligned} E[\mathbf{y}^p | \mathbf{y} = 22,000] &= \mu + \frac{\tau_y^2}{\tau_y^2 + \sigma_y^2} (22,000 - \mu) \\ &= \frac{\sigma_y^2}{\tau_y^2 + \sigma_y^2} \mu + \frac{\tau_y^2}{\tau_y^2 + \sigma_y^2} 22,000 \\ &= (0.2)(12,000) + (0.8)(22,000) = 20,000. \end{aligned}$$

□

• c. 3 points To make things more concrete, assume the parameters are

$$(8.0.23) \quad \beta = 0.7$$

$$(8.0.24) \quad \sigma_y = 2,000$$

$$(8.0.25) \quad \sigma_c = 1,000$$

$$(8.0.26) \quad \mu = 12,000$$

$$(8.0.27) \quad \tau_y = 4,000.$$

If a family's income is  $y = 22,000$ , what is your best guess of this family's permanent income  $\mathbf{y}^p$ ? Give an intuitive explanation why this best guess is smaller than 22,000.

ANSWER. Since the observed income of 22,000 is above the average of 12,000, chances greater that it is someone with a positive transitory income than someone with a negative one.

• d. 2 points If a family's income is  $\mathbf{y}$ , show that your best guess about the family's consumption is

$$(8.0.28) \quad \mathbf{c}^* = \beta \left( \frac{\sigma_y^2}{\tau_y^2 + \sigma_y^2} \mu + \frac{\tau_y^2}{\tau_y^2 + \sigma_y^2} \mathbf{y} \right).$$

Instead of an exact mathematical proof you may also reason out how it can be obtained from (8.0.22). Give the numbers for a family whose actual income is 22,000.

ANSWER. This is 0.7 times the best guess about the family's permanent income, since transitory consumption is uncorrelated with everything else and therefore must be predicted by  $\beta$ . This is an acceptable answer, but one can also derive it from scratch:

(8.0.29)

$$(8.0.30) \quad \begin{aligned} E[\mathbf{c} | \mathbf{y} = 22,000] &= E[\mathbf{c}] + \frac{\text{cov}[\mathbf{c}, \mathbf{y}]}{\text{var}[\mathbf{y}]} (22,000 - E[\mathbf{y}]) \\ &= \beta\mu + \frac{\beta\tau_y^2}{\tau_y^2 + \sigma_y^2} (22,000 - \mu) = 8,400 + 0.7 \frac{16,000,000}{20,000,000} (22,000 - 12,000) = 14,000 \end{aligned}$$

$$(8.0.31) \quad \text{or} \quad = \beta \left( \frac{\sigma_y^2}{\tau_y^2 + \sigma_y^2} \mu + \frac{\tau_y^2}{\tau_y^2 + \sigma_y^2} 22,000 \right)$$

$$(8.0.32) \quad = 0.7 \left( (0.2)(12,000) + (0.8)(22,000) \right) = (0.7)(20,000) = 14,000.$$

The remainder of this Problem uses material that comes later in these Notes

• e. 4 points From now on we will assume that the true values of the parameters are not known, but two vectors  $\mathbf{y}$  and  $\mathbf{c}$  of independent observations are available. We will show that it is not correct in this situation to estimate  $\beta$  by regressing  $\mathbf{c}$  on  $\mathbf{y}$  with the intercept suppressed. This would give the estimator

$$(8.0.33) \quad \hat{\beta} = \frac{\sum \mathbf{c}_i \mathbf{y}_i}{\sum \mathbf{y}_i^2}$$

Show that the plim of this estimator is

$$(8.0.34) \quad \text{plim}[\hat{\beta}] = \frac{E[\mathbf{c}\mathbf{y}]}{E[\mathbf{y}^2]}$$

Which theorems do you need for this proof? Show that  $\hat{\beta}$  is an inconsistent estimator of  $\beta$ , which yields too small values for  $\beta$ .

ANSWER. First rewrite the formula for  $\hat{\beta}$  in such a way that numerator and denominator each has a plim: by the weak law of large numbers the plim of the average is the expected value, therefore we have to divide both numerator and denominator by  $n$ . Then we can use the Slutsky theorem that the plim of the fraction is the fraction of the plims.

$$\hat{\beta} = \frac{\frac{1}{n} \sum c_i y_i}{\frac{1}{n} \sum y_i^2}; \quad \text{plim}[\hat{\beta}] = \frac{E[\mathbf{c}\mathbf{y}]}{E[\mathbf{y}^2]} = \frac{E[\mathbf{c}]E[\mathbf{y}] + \text{cov}[\mathbf{c}, \mathbf{y}]}{(E[\mathbf{y}])^2 + \text{var}[\mathbf{y}]} = \frac{\mu\beta\mu + \beta\tau_y^2}{\mu^2 + \tau_y^2 + \sigma_y^2} = \beta \frac{\mu^2 + \tau_y^2}{\mu^2 + \tau_y^2 + \sigma_y^2}.$$

□

• f. 4 points Give the formulas of the method of moments estimators of the five parameters of this model:  $\mu$ ,  $\beta$ ,  $\tau_y^2$ ,  $\sigma_y^2$ , and  $\sigma_p^2$ . (For this you have to express these five parameters in terms of the five moments  $E[\mathbf{y}]$ ,  $E[\mathbf{c}]$ ,  $\text{var}[\mathbf{y}]$ ,  $\text{var}[\mathbf{c}]$ , and  $\text{cov}[\mathbf{y}, \mathbf{c}]$ , and then simply replace the population moments by the sample moments.) Are these consistent estimators?

ANSWER. From (8.0.21) follows  $E[\mathbf{c}] = \beta E[\mathbf{y}]$ , therefore  $\beta = \frac{E[\mathbf{c}]}{E[\mathbf{y}]}$ . This together with  $\text{cov}[\mathbf{y}, \mathbf{c}] = \beta\tau_y^2$  gives  $\tau_y^2 = \frac{\text{cov}[\mathbf{y}, \mathbf{c}]}{\beta} = \frac{\text{cov}[\mathbf{y}, \mathbf{c}]E[\mathbf{y}]}{E[\mathbf{c}]}$ . This together with  $\text{var}[\mathbf{y}] = \tau_y^2 + \sigma_y^2$  gives  $\sigma_y^2 = \text{var}[\mathbf{y}] - \tau_y^2 = \text{var}[\mathbf{y}] - \frac{\text{cov}[\mathbf{y}, \mathbf{c}]E[\mathbf{y}]}{E[\mathbf{c}]}$ . And from the last equation  $\text{var}[\mathbf{c}] = \beta^2\tau_y^2 + \sigma_c^2$  one get  $\sigma_c^2 = \text{var}[\mathbf{c}] - \frac{\text{cov}[\mathbf{y}, \mathbf{c}]E[\mathbf{c}]}{E[\mathbf{y}]}$ . All these are consistent estimators, as long as  $E[\mathbf{y}] \neq 0$  and  $\beta \neq 0$ . □

• g. 4 points Now assume you are not interested in estimating  $\beta$  itself, but in addition to the two  $n$ -vectors  $\mathbf{y}$  and  $\mathbf{c}$  you have an observation of  $\mathbf{y}_{n+1}$  and you want to predict the corresponding  $\mathbf{c}_{n+1}$ . One obvious way to do this would be to plug the method-of-moments estimators of the unknown parameters into formula (8.0.28) for the best linear predictor. Show that this is equivalent to using the ordinary least squares predictor  $\mathbf{c}^* = \hat{\alpha} + \hat{\beta}\mathbf{y}_{n+1}$  where  $\hat{\alpha}$  and  $\hat{\beta}$  are intercept and slope in the simple regression of  $\mathbf{c}$  on  $\mathbf{y}$ , i.e.,

$$(8.0.35) \quad \hat{\beta} = \frac{\sum(\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{c}_i - \bar{\mathbf{c}})}{\sum(\mathbf{y}_i - \bar{\mathbf{y}})^2}$$

$$(8.0.36) \quad \hat{\alpha} = \bar{\mathbf{c}} - \hat{\beta}\bar{\mathbf{y}}$$

Note that we are regressing  $\mathbf{c}$  on  $\mathbf{y}$  with an intercept, although the original model does not have an intercept.

ANSWER. Here I am writing population moments where I should be writing sample moments. First substitute the method of moments estimators in the denominator in (8.0.28):  $\tau_y^2 + \sigma_y^2 = \text{var}[\mathbf{y}]$ . Therefore the first summand becomes

$$\beta\sigma_y^2\mu \frac{1}{\text{var}[\mathbf{y}]} = \frac{E[\mathbf{c}]}{E[\mathbf{y}]} \left( \text{var}[\mathbf{y}] - \frac{\text{cov}[\mathbf{y}, \mathbf{c}]E[\mathbf{y}]}{E[\mathbf{c}]} \right) \frac{1}{\text{var}[\mathbf{y}]} = E[\mathbf{c}] \left( 1 - \frac{\text{cov}[\mathbf{y}, \mathbf{c}]E[\mathbf{y}]}{\text{var}[\mathbf{y}]E[\mathbf{c}]} \right) = E[\mathbf{c}] - \frac{\text{cov}[\mathbf{y}, \mathbf{c}]}{\text{var}[\mathbf{y}]}$$

But since  $\frac{\text{cov}[\mathbf{y}, \mathbf{c}]}{\text{var}[\mathbf{y}]} = \hat{\beta}$  and  $\hat{\alpha} + \hat{\beta}E[\mathbf{y}] = E[\mathbf{c}]$  this expression is simply  $\hat{\alpha}$ . The second term is easy to show:

$$\beta \frac{\tau_y^2}{\text{var}[\mathbf{y}]} \mathbf{y} = \frac{\text{cov}[\mathbf{y}, \mathbf{c}]}{\text{var}[\mathbf{y}]} \mathbf{y} = \hat{\beta}\mathbf{y}$$

• h. 2 points What is the “Iron Law of Econometrics,” and how does the above relate to it?

ANSWER. The Iron Law says that all effects are underestimated because of errors in the independent variable. Friedman says Keynesians obtain their low marginal propensity to consume due to the “Iron Law of Econometrics”: they ignore that actual income is a measurement with error of the true underlying variable, permanent income.

PROBLEM 154. This question follows the original article [SW76] much more closely than [HVdP02] does. Sargent and Wallace first reproduce the usual argument why “activist” policy rules, in which the Fed “looks at many things” and “leaves the money supply against the wind,” are superior to policy rules without feedback as promoted by monetarists.

They work with a very stylized model in which national income is represented by the following time series:

$$(8.0.37) \quad \mathbf{y}_t = \alpha + \lambda\mathbf{y}_{t-1} + \beta\mathbf{m}_t + \mathbf{u}_t$$

Here  $\mathbf{y}_t$  is GNP, measured as its deviation from “potential” GNP or as unemployment rate, and  $\mathbf{m}_t$  is the rate of growth of the money supply. The random disturbance  $\mathbf{u}_t$  is assumed independent of  $\mathbf{y}_{t-1}$ , it has zero expected value, and its variance  $\text{var}[\mathbf{u}_t]$  is constant over time, we will call it  $\text{var}[\mathbf{u}]$  (no time subscript).

• a. 4 points First assume that the Fed tries to maintain a constant money supply, i.e.,  $\mathbf{m}_t = g_0 + \varepsilon_t$  where  $g_0$  is a constant, and  $\varepsilon_t$  is a random disturbance since the Fed does not have full control over the money supply. The  $\varepsilon_t$  have zero expected value; they are serially uncorrelated, and they are independent of the  $\mathbf{y}_t$ . This constant money supply rule does not necessarily make  $\mathbf{y}_t$  a stationary time series (i.e., a time series where mean, variance, and covariances do not depend on  $t$ ), but if  $|\lambda| < 1$  then  $\mathbf{y}_t$  converges towards a stationary time series, i.e., any initial deviations from the “steady state” die out over time. You are not required here to prove that the time series converges towards a stationary time series, but you are asked to compute  $E[\mathbf{y}_t]$  in this stationary time series.

• b. 8 points Now assume the policy makers want to steer the economy towards a desired steady state, call it  $y^*$ , which they think makes the best tradeoff between unemployment and inflation, by setting  $m_t$  according to a rule with feedback:

$$(8.0.38) \quad m_t = g_0 + g_1 y_{t-1} + \varepsilon_t$$

Show that the following values of  $g_0$  and  $g_1$

$$(8.0.39) \quad g_0 = (y^* - \alpha)/\beta \quad g_1 = -\lambda/\beta$$

represent an optimal monetary policy, since they bring the expected value of the steady state  $E[y_t]$  to  $y^*$  and minimize the steady state variance  $\text{var}[y_t]$ .

• c. 3 points This is the conventional reasoning which comes to the result that a policy rule with feedback, i.e., a policy rule in which  $g_1 \neq 0$ , is better than a policy rule without feedback. Sargent and Wallace argue that there is a flaw in this reasoning. Which flaw?

• d. 5 points A possible system of structural equations from which (8.0.37) can be derived are equations (8.0.40)–(8.0.42) below. Equation (8.0.40) indicates that unanticipated increases in the growth rate of the money supply increase output, while anticipated ones do not. This is a typical assumption of the rational expectations school (Lucas supply curve).

$$(8.0.40) \quad y_t = \xi_0 + \xi_1(m_t - E_{t-1} m_t) + \xi_2 y_{t-1} + u_t$$

The Fed uses the policy rule

$$(8.0.41) \quad m_t = g_0 + g_1 y_{t-1} + \varepsilon_t$$

and the agents know this policy rule, therefore

$$(8.0.42) \quad E_{t-1} m_t = g_0 + g_1 y_{t-1}.$$

Show that in this system, the parameters  $g_0$  and  $g_1$  have no influence on the time path of  $y$ .

• e. 4 points On the other hand, the econometric estimations which the policy makers are running seem to show that these coefficients have an impact. During a certain period during which a constant policy rule  $g_0, g_1$  is followed, the econometricians regress  $y_t$  on  $y_{t-1}$  and  $m_t$  in order to estimate the coefficients in (8.0.37). Which values of  $\alpha, \lambda$ , and  $\beta$  will such a regression yield?

## CHAPTER 9

## A Simple Example of Estimation

We will discuss here a simple estimation problem, which can be considered the prototype of all least squares estimation. Assume we have  $n$  independent observations  $y_1, \dots, y_n$  of a Normally distributed random variable  $\mathbf{y} \sim N(\mu, \sigma^2)$  with unknown location parameter  $\mu$  and dispersion parameter  $\sigma^2$ . Our goal is to estimate the *location parameter* and also estimate some measure of the precision of this estimator.

## 9.1. Sample Mean as Estimator of the Location Parameter

The obvious (and in many cases also the best) estimate of the location parameter of a distribution is the sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . Why is this a reasonable estimate?

1. The location parameter of the Normal distribution is its expected value, and by the weak law of large numbers, the probability limit for  $n \rightarrow \infty$  of the sample mean is the expected value.

2. The expected value  $\mu$  is sometimes called the “population mean,” while  $\bar{y}$  is the sample mean. This terminology indicates that there is a correspondence between population quantities and sample quantities, which is often used for estimation. This is the principle of estimating the unknown distribution of the population by the empirical distribution of the sample. Compare Problem 63.

3. This estimator is also unbiased. By definition, an estimator  $t$  of the parameter  $\theta$  is unbiased if  $E[t] = \theta$ .  $\bar{y}$  is an unbiased estimator of  $\mu$ , since  $E[\bar{y}] = \mu$ .

4. Given  $n$  observations  $y_1, \dots, y_n$ , the sample mean is the number  $a = \bar{y}$  which minimizes  $(y_1 - a)^2 + (y_2 - a)^2 + \dots + (y_n - a)^2$ . One can say it is the number whose squared distance to the given sample numbers is smallest. This idea is generalized in the least squares principle of estimation. It follows from the following frequently used fact:

5. In the case of normality the sample mean is also the maximum likelihood estimate.

PROBLEM 155. 4 points Let  $y_1, \dots, y_n$  be an arbitrary vector and  $\alpha$  an arbitrary number. As usual,  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . Show that

$$(9.1.1) \quad \sum_{i=1}^n (y_i - \alpha)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \alpha)^2$$

ANSWER.

$$(9.1.2) \quad \sum_{i=1}^n (y_i - \alpha)^2 = \sum_{i=1}^n ((y_i - \bar{y}) + (\bar{y} - \alpha))^2$$

$$(9.1.3) \quad = \sum_{i=1}^n (y_i - \bar{y})^2 + 2 \sum_{i=1}^n ((y_i - \bar{y})(\bar{y} - \alpha)) + \sum_{i=1}^n (\bar{y} - \alpha)^2$$

$$(9.1.4) \quad = \sum_{i=1}^n (y_i - \bar{y})^2 + 2(\bar{y} - \alpha) \sum_{i=1}^n (y_i - \bar{y}) + n(\bar{y} - \alpha)^2$$

Since the middle term is zero, (9.1.1) follows.

PROBLEM 156. 2 points Let  $\mathbf{y}$  be a  $n$ -vector. (It may be a vector of observations of a random variable  $\mathbf{y}$ , but it does not matter how the  $y_i$  were obtained.) Prove that the scalar  $\alpha$  which minimizes the sum

$$(9.1.5) \quad (y_1 - \alpha)^2 + (y_2 - \alpha)^2 + \dots + (y_n - \alpha)^2 = \sum (y_i - \alpha)^2$$

is the arithmetic mean  $\alpha = \bar{y}$ .

ANSWER. Use (9.1.1).

PROBLEM 157. Give an example of a distribution in which the sample mean is not a good estimate of the location parameter. Which other estimate (or estimates) would be preferable in that situation?

## 9.2. Intuition of the Maximum Likelihood Estimator

In order to make intuitively clear what is involved in maximum likelihood estimation, look at the simplest case  $\mathbf{y} = \mu + \varepsilon$ ,  $\varepsilon \sim N(0, 1)$ , where  $\mu$  is an unknown parameter. In other words: we know that one of the functions shown in Figure 1 is the density function of  $\mathbf{y}$ , but we do not know which:

Assume we have only one observation  $y$ . What is then the MLE of  $\mu$ ? It is the  $\tilde{\mu}$  for which the value of the likelihood function, evaluated at  $y$ , is greatest. I.e., you look at all possible density functions and pick the one which is highest at point  $y$  and use the  $\mu$  which belongs to this density as your estimate.

2) Now assume two independent observations of  $\mathbf{y}$  are given,  $y_1$  and  $y_2$ . The family of density functions is still the same. Which of these density functions do you choose now? The one for which the *product* of the ordinates over  $y_1$  and  $y_2$  gives

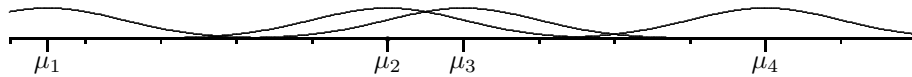


FIGURE 1. Possible Density Functions for  $y$

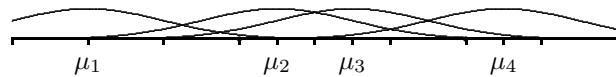


FIGURE 2. Two observations,  $\sigma^2 = 1$

FIGURE 3. Two observations,  $\sigma^2$  unknown

the highest value. For this the peak of the density function must be exactly in the middle between the two observations.

3) Assume again that we made two independent observations  $y_1$  and  $y_2$  of  $y$ , but this time not only the expected value but also the variance of  $y$  is unknown, call it  $\sigma^2$ . This gives a larger family of density functions to choose from: they do not only differ by location, but some are low and fat and others tall and skinny.

For which density function is the product of the ordinates over  $y_1$  and  $y_2$  the largest again? Before even knowing our estimate of  $\sigma^2$  we can already tell what  $\tilde{\mu}$  is: it must again be  $(y_1 + y_2)/2$ . Then among those density functions which are centered over  $(y_1 + y_2)/2$ , there is one which is highest over  $y_1$  and  $y_2$ . Figure 4 shows the densities for standard deviations 0.01, 0.05, 0.1, 0.5, 1, and 5. All curves, except the last one, are truncated at the point where the resolution of  $\text{T}_{\text{E}}\text{X}$  can no longer distinguish between their level and zero. For the last curve this point would only be reached at the coordinates  $\pm 25$ .

4) If we have many observations, then the density pattern of the observations, as indicated by the histogram below, approximates the actual density function of  $y$  itself. That likelihood function must be chosen which has a high value where the points are dense, and which has a low value where the points are not so dense.

**9.2.1. Precision of the Estimator.** How good is  $\bar{y}$  as estimate of  $\mu$ ? To answer this question we need some criterion how to measure “goodness.” Assume your

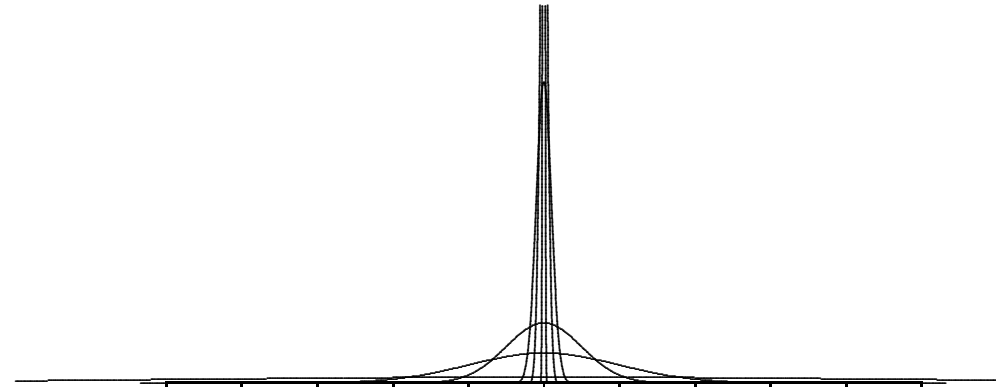


FIGURE 4. Only those centered over the two observations need to be considered

FIGURE 5. Many Observations

business depends on the precision of the estimate  $\hat{\mu}$  of  $\mu$ . It incurs a penalty (expected cost) amounting to  $(\hat{\mu} - \mu)^2$ . You don’t know what this error will be beforehand, but the expected value of this “loss function” may be an indication how good the estimate is. Generally, the expected value of a loss function is called the “risk,” and for the quadratic loss function  $E[(\hat{\mu} - \mu)^2]$  it has the name “mean squared error of  $\hat{\mu}$  as an estimate of  $\mu$ ,” write it  $\text{MSE}[\hat{\mu}; \mu]$ . What is the mean squared error of  $\bar{y}$  as an estimate of  $\mu$ , it is  $E[(\bar{y} - E[\bar{y}])^2] = \text{var}[\bar{y}] = \frac{\sigma^2}{n}$ .

Note that the MSE of  $\bar{y}$  as an estimate of  $\mu$  does not depend on  $\mu$ . This is convenient, since usually the MSE depends on unknown parameters, and therefore one usually does not know how good the estimator is. But it has more important advantages. For any estimator  $\tilde{y}$  of  $\mu$  follows  $\text{MSE}[\tilde{y}; \mu] = \text{var}[\tilde{y}] + (E[\tilde{y}] - \mu)^2$ . If  $\tilde{y}$  is linear (perhaps with a constant term), then  $\text{var}[\tilde{y}]$  is a constant which does not depend on  $\mu$ , therefore the MSE is a constant if  $\tilde{y}$  is unbiased and a quadratic function of  $\mu$  (parabola) if  $\tilde{y}$  is biased. Since a parabola is an unbounded function, a biased linear estimator has therefore the disadvantage that for certain values of  $\mu$  its MSE may be very high. Some estimators are very good when  $\mu$  is in one area and very bad when  $\mu$  is in another area. Since our unbiased estimator  $\bar{y}$  has bounded MSE, it will not let us down, wherever nature has hidden the  $\mu$ .

On the other hand, the MSE does depend on the unknown  $\sigma^2$ . So we have to estimate  $\sigma^2$ .

### 9.3. Variance Estimation and Degrees of Freedom

It is not so clear what the best estimator of  $\sigma^2$  is. At least two possibilities are in common use:

$$(9.3.1) \quad s_m^2 = \frac{1}{n} \sum (y_i - \bar{y})^2$$

or

$$(9.3.2) \quad s_u^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2.$$

Let us compute the expected value of our two estimators. Equation (9.1.1) with  $\alpha = E[y]$  allows us to simplify the sum of squared errors so that it becomes easy to take expected values:

$$(9.3.3) \quad E\left[\sum_{i=1}^n (y_i - \bar{y})^2\right] = \sum_{i=1}^n E[(y_i - \mu)^2] - n E[(\bar{y} - \mu)^2]$$

$$(9.3.4) \quad = \sum_{i=1}^n \sigma^2 - n \frac{\sigma^2}{n} = (n-1)\sigma^2.$$

because  $E[(y_i - \mu)^2] = \text{var}[y_i] = \sigma^2$  and  $E[(\bar{y} - \mu)^2] = \text{var}[\bar{y}] = \frac{\sigma^2}{n}$ . Therefore, if we use as estimator of  $\sigma^2$  the quantity

$$(9.3.5) \quad s_u^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

then this is an unbiased estimate.

PROBLEM 158. 4 points Show that

$$(9.3.6) \quad s_u^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

is an unbiased estimator of the variance. List the assumptions which have to be made about  $y_i$  so that this proof goes through. Do you need Normality of the individual observations  $y_i$  to prove this?

ANSWER. Use equation (9.1.1) with  $\alpha = E[y]$ :

$$(9.3.7) \quad E\left[\sum_{i=1}^n (y_i - \bar{y})^2\right] = \sum_{i=1}^n E[(y_i - \mu)^2] - n E[(\bar{y} - \mu)^2]$$

$$(9.3.8) \quad = \sum_{i=1}^n \sigma^2 - n \frac{\sigma^2}{n} = (n-1)\sigma^2.$$

You do not need Normality for this. □

For testing, confidence intervals, etc., one also needs to know the probability distribution of  $s_u^2$ . For this look up once more Section 4.9 about the Chi-Square distribution. There we introduced the terminology that a random variable  $q$  is distributed as a  $\sigma^2 \chi^2$  iff  $q/\sigma^2$  is a  $\chi^2$ . In our model with  $n$  independent normal variables  $y_i$  with same mean and variance, the variable  $\sum (y_i - \bar{y})^2$  is a  $\sigma^2 \chi_{n-1}^2$ . Problem 159 gives a proof of this in the simplest case  $n = 2$ , and Problem 160 looks at the case  $n = 3$ . But it is valid for higher  $n$  too. Therefore  $s_u^2$  is a  $\frac{\sigma^2}{n-1} \chi_{n-1}^2$ . This is remarkable: the distribution of  $s_u^2$  does not depend on  $\mu$ . Now use (4.9.5) to find the variance of  $s_u^2$ : it is  $\frac{2\sigma^4}{n-1}$ .

PROBLEM 159. Let  $y_1$  and  $y_2$  be two independent Normally distributed variables with mean  $\mu$  and variance  $\sigma^2$ , and let  $\bar{y}$  be their arithmetic mean.

• a. 2 points Show that

$$(9.3.9) \quad \text{SSE} = \sum_{i=1}^2 (y_i - \bar{y})^2 \sim \sigma^2 \chi_1^2$$

Hint: Find a Normally distributed random variable  $z$  with expected value 0 and variance 1 such that  $\text{SSE} = \sigma^2 z^2$ .

ANSWER.

$$(9.3.10) \quad \bar{y} = \frac{y_1 + y_2}{2}$$

$$(9.3.11) \quad y_1 - \bar{y} = \frac{y_1 - y_2}{2}$$

$$(9.3.12) \quad y_2 - \bar{y} = -\frac{y_1 - y_2}{2}$$

$$(9.3.13) \quad (y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 = \frac{(y_1 - y_2)^2}{4} + \frac{(y_1 - y_2)^2}{4}$$

$$(9.3.14) \quad = \frac{(y_1 - y_2)^2}{2} = \sigma^2 \left( \frac{y_1 - y_2}{\sqrt{2}\sigma} \right)^2,$$

and since  $z = (y_1 - y_2)/\sqrt{2}\sigma \sim N(0, 1)$ , its square is a  $\chi_1^2$ .

• b. 4 points Write down the covariance matrix of the vector

$$(9.3.15) \quad \begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \end{bmatrix}$$

and show that it is singular.

ANSWER. (9.3.11) and (9.3.12) give

$$(9.3.16) \quad \begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = D\mathbf{y}$$

and  $\mathcal{V}[\mathbf{D}\mathbf{y}] = \mathbf{D}\mathcal{V}[\mathbf{y}]\mathbf{D}^\top = \sigma^2\mathbf{D}$  because  $\mathcal{V}[\mathbf{y}] = \sigma^2\mathbf{I}$  and  $\mathbf{D} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}$  is symmetric and idempotent.  $\mathbf{D}$  is singular because its determinant is zero.  $\square$

• c. 1 point The joint distribution of  $\mathbf{y}_1$  and  $\mathbf{y}_2$  is bivariate normal, why did we then get a  $\chi^2$  with one, instead of two, degrees of freedom?

ANSWER. Because  $\mathbf{y}_1 - \bar{y}$  and  $\mathbf{y}_2 - \bar{y}$  are not independent; one is exactly the negative of the other; therefore summing their squares is really only the square of one univariate normal.  $\square$

PROBLEM 160. Assume  $\mathbf{y}_1, \mathbf{y}_2,$  and  $\mathbf{y}_3$  are independent  $\mathcal{N}(\mu, \sigma^2)$ . Define three new variables  $\mathbf{z}_1, \mathbf{z}_2,$  and  $\mathbf{z}_3$  as follows:  $\mathbf{z}_1$  is that multiple of  $\bar{y}$  which has variance  $\sigma^2$ .  $\mathbf{z}_2$  is that linear combination of  $\mathbf{z}_1$  and  $\mathbf{y}_2$  which has zero covariance with  $\mathbf{z}_1$  and has variance  $\sigma^2$ .  $\mathbf{z}_3$  is that linear combination of  $\mathbf{z}_1, \mathbf{z}_2,$  and  $\mathbf{y}_3$  which has zero covariance with both  $\mathbf{z}_1$  and  $\mathbf{z}_2$  and has again variance  $\sigma^2$ . These properties define  $\mathbf{z}_1, \mathbf{z}_2,$  and  $\mathbf{z}_3$  uniquely up factors  $\pm 1$ , i.e., if  $\mathbf{z}_1$  satisfies the above conditions, then  $-\mathbf{z}_1$  does too, and these are the only two solutions.

• a. 2 points Write  $\mathbf{z}_1$  and  $\mathbf{z}_2$  (not yet  $\mathbf{z}_3$ ) as linear combinations of  $\mathbf{y}_1, \mathbf{y}_2,$  and  $\mathbf{y}_3$ .

• b. 1 point To make the computation of  $\mathbf{z}_3$  less tedious, first show the following: if  $\mathbf{z}_3$  has zero covariance with  $\mathbf{z}_1$  and  $\mathbf{z}_2$ , it also has zero covariance with  $\mathbf{y}_2$ .

• c. 1 point Therefore  $\mathbf{z}_3$  is a linear combination of  $\mathbf{y}_1$  and  $\mathbf{y}_3$  only. Compute its coefficients.

• d. 1 point How does the joint distribution of  $\mathbf{z}_1, \mathbf{z}_2,$  and  $\mathbf{z}_3$  differ from that of  $\mathbf{y}_1, \mathbf{y}_2,$  and  $\mathbf{y}_3$ ? Since they are jointly normal, you merely have to look at the expected values, variances, and covariances.

• e. 2 points Show that  $\mathbf{z}_1^2 + \mathbf{z}_2^2 + \mathbf{z}_3^2 = \mathbf{y}_1^2 + \mathbf{y}_2^2 + \mathbf{y}_3^2$ . Is this a surprise?

• f. 1 point Show further that  $s_u^2 = \frac{1}{2} \sum_{i=1}^3 (\mathbf{y}_i - \bar{y})^2 = \frac{1}{2} (\mathbf{z}_2^2 + \mathbf{z}_3^2)$ . (There is a simple trick!) Conclude from this that  $s_u^2 \sim \frac{\sigma^2}{2} \chi_2^2$ , independent of  $\bar{y}$ .

For a matrix-interpretation of what is happening, see equation (7.4.9) together with Problem 161.

PROBLEM 161. 3 points Verify that the matrix  $\mathbf{D} = \mathbf{I} - \frac{1}{n}\mathbf{u}\mathbf{u}^\top$  is symmetric and idempotent, and that the sample covariance of two vectors of observations  $\mathbf{x}$  and  $\mathbf{y}$  can be written in matrix notation as

$$(9.3.17) \quad \text{sample covariance}(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \mathbf{x}^\top \mathbf{D} \mathbf{y}$$

In general, one can always find  $n - 1$  normal variables with variance  $\sigma^2$ , independent of each other and of  $\bar{y}$ , whose sum of squares is equal to  $\sum (\mathbf{y}_i - \bar{y})^2$ . Simply

start with  $\bar{y}\sqrt{n}$  and generate  $n - 1$  linear combinations of the  $\mathbf{y}_i$  which are pairwise uncorrelated and have variances  $\sigma^2$ . You are simply building an orthonormal coordinate system with  $\bar{y}\sqrt{n}$  as its first vector; there are many different ways to do this.

Next let us show that  $\bar{y}$  and  $s_u^2$  are statistically independent. This is an advantage. Assume, hypothetically,  $\bar{y}$  and  $s_u^2$  were negatively correlated. Then, if the observed value of  $\bar{y}$  is too high, chances are that the one of  $s_u^2$  is too low, and a low value at  $s_u^2$  will not reveal how far off the mark  $\bar{y}$  may be. To prove independence, we will first show that  $\bar{y}$  and  $\mathbf{y}_i - \bar{y}$  are uncorrelated:

$$(9.3.18) \quad \text{cov}[\bar{y}, \mathbf{y}_i - \bar{y}] = \text{cov}[\bar{y}, \mathbf{y}_i] - \text{var}[\bar{y}]$$

$$(9.3.19) \quad = \text{cov}\left[\frac{1}{n}(\mathbf{y}_1 + \cdots + \mathbf{y}_i + \cdots + \mathbf{y}_n), \mathbf{y}_i\right] - \frac{\sigma^2}{n} = 0$$

By normality,  $\bar{y}$  is therefore independent of  $\mathbf{y}_i - \bar{y}$  for all  $i$ . Since all variables involved are jointly normal, it follows from this that  $\bar{y}$  is independent of the vector  $[\mathbf{y}_1 - \bar{y} \quad \cdots \quad \mathbf{y}_n - \bar{y}]^\top$ ; therefore it is also independent of any function of this vector, such as  $s_u^2$ .

The above calculations explain why the parameter of the  $\chi^2$  distribution has the colorful name “degrees of freedom.” This term is sometimes used in a very broad sense, referring to estimation in general, and sometimes in a narrower sense in conjunction with the linear model. Here is first an interpretation of the general use of the term. A “statistic” is defined to be a function of the observations and of other known parameters of the problem, but not of the unknown parameters. Estimates are statistics. If one has  $n$  observations, then one can find at most  $n$  mathematically independent statistics; any other statistic is then a function of these  $n$ . If therefore a model has  $k$  independent unknown parameters, then one must have at least  $n - k$  observations to be able to estimate all parameters of the model. The number  $n - k$ , i.e., the number of observations not “used up” for estimation, is called the number of “degrees of freedom.”

There are at least three reasons why one does not want to make the model such that it uses up too many degrees of freedom. (1) the estimators become too inaccurate if one does; (2) if there are no degrees of freedom left, it is no longer possible to make any “diagnostic” tests whether the model really fits the data, because it always gives a perfect fit whatever the given set of data; (3) if there are no degrees of freedom left then one can usually also no longer make estimates of the precision of the estimators.

Specifically in our linear estimation problem, the number of degrees of freedom is  $n - 1$ , since one observation has been used up for estimating the mean. If one runs a regression, the number of degrees of freedom is  $n - k$ , where  $k$  is the number of regression coefficients. In the linear model, the number of degrees of freedom becomes immediately relevant for the estimation of  $\sigma^2$ . If  $k$  observations are used

up for estimating the slope parameters, then the other  $n - k$  observations can be combined into a  $n - k$ -variate Normal whose expected value does not depend on the slope parameter at all but is zero, which allows one to estimate the variance.

If we assume that the original observations are normally distributed, i.e.,  $y_i \sim \text{NID}(\mu, \sigma^2)$ , then we know that  $s_u^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$ . Therefore  $E[s_u^2] = \sigma^2$  and  $\text{var}[s_u^2] = 2\sigma^4/(n - 1)$ . This estimate of  $\sigma^2$  therefore not only gives us an estimate of the precision of  $\bar{y}$ , but it has an estimate of its own precision built in.

Interestingly, the MSE of the alternative estimator  $s_m^2 = \frac{\sum(y_i - \bar{y})^2}{n}$  is smaller than that of  $s_u^2$ , although  $s_m^2$  is a biased estimator and  $s_u^2$  an unbiased estimator of  $\sigma^2$ . For every estimator  $t$ ,  $\text{MSE}[t; \theta] = \text{var}[t] + (E[t - \theta])^2$ , i.e., it is variance plus squared bias. The MSE of  $s_u^2$  is therefore equal to its variance, which is  $\frac{2\sigma^4}{n-1}$ . The alternative  $s_m^2 = \frac{n-1}{n}s_u^2$  has bias  $-\frac{\sigma^2}{n}$  and variance  $\frac{2\sigma^4(n-1)}{n^2}$ . Its MSE is  $\frac{(2-1/n)\sigma^4}{n}$ . Comparing that with the formula for the MSE of  $s_u^2$  one sees that the numerator is smaller and the denominator is bigger, therefore  $s_m^2$  has smaller MSE.

**PROBLEM 162.** 4 points Assume  $y_i \sim \text{NID}(\mu, \sigma^2)$ . Show that the so-called Theil Schweitzer estimator [TS61]

$$(9.3.20) \quad s_t^2 = \frac{1}{n+1} \sum (y_i - \bar{y})^2$$

has even smaller MSE than  $s_u^2$  and  $s_m^2$  as an estimator of  $\sigma^2$ .

**ANSWER.**  $s_t^2 = \frac{n-1}{n+1}s_u^2$ ; therefore its bias is  $-\frac{2\sigma^2}{n+1}$  and its variance is  $\frac{2(n-1)\sigma^4}{(n+1)^2}$ , and the MSE is  $\frac{2\sigma^4}{n+1}$ . That this is smaller than the MSE of  $s_m^2$  means  $\frac{2n-1}{n^2} \geq \frac{2}{n+1}$ , which follows from  $(2n - 1)(n + 1) = 2n^2 + n - 1 > 2n^2$  for  $n > 1$ .  $\square$

**PROBLEM 163.** 3 points Computer assignment: Given 20 independent observations of a random variable  $y \sim \text{N}(\mu, \sigma^2)$ . Assume you know that  $\sigma^2 = 2$ . Plot the density function of  $s_u^2$ . Hint: In R, the command `dchisq(x, df=25)` returns the density of a Chi-square distribution with 25 degrees of freedom evaluated at  $x$ . But the number 25 was only taken as an example, this is not the number of degrees of freedom you need here. You also do not need the density of a Chi-Square but that of a certain multiple of a Chi-square. (Use the transformation theorem for density functions!)

**ANSWER.**  $s_u^2 \sim \frac{2}{19} \chi_{19}^2$ . To express the density of the variable whose density is known by that whose density one wants to know, say  $\frac{19}{2}s_u^2 \sim \chi_{19}^2$ . Therefore

$$(9.3.21) \quad f_{s_u^2}(x) = \frac{19}{2} f_{\chi_{19}^2} \left( \frac{19}{2} x \right).$$

$\square$

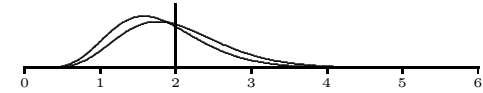


FIGURE 6. Densities of Unbiased and Theil Schweitzer Estimators

• a. 2 points In the same plot, plot the density function of the Theil-Schweitzer estimate  $s_t^2$  defined in equation (9.3.20). This gives a plot as in Figure 6. Can one from the comparison of these density functions that the Theil-Schweitzer estimator has a better MSE?

**ANSWER.** Start with plotting the Theil-Schweitzer plot, because it is higher, and therefore will give the right dimensions of the plot. You can run this by giving the command `ecmetscript` (The two areas between the densities have equal size, but the area where the Theil-Schweitzer density is higher is overall closer to the true value than the area where the unbiased density is higher.

**PROBLEM 164.** 4 points The following problem illustrates the general fact that if one starts with an unbiased estimator and “shrinks” it a little, one will end up with a better MSE. Assume  $E[y] = \mu$ ,  $\text{var}(y) = \sigma^2$ , and you make  $n$  independent observations  $y_i$ . The best linear unbiased estimator of  $\mu$  on the basis of the observations is the sample mean  $\bar{y}$ . Show that, whenever  $\alpha$  satisfies

$$(9.3.22) \quad \frac{n\mu^2 - \sigma^2}{n\mu^2 + \sigma^2} < \alpha < 1$$

then  $\text{MSE}[\alpha\bar{y}; \mu] < \text{MSE}[\bar{y}; \mu]$ . Unfortunately, this condition depends on  $\mu$  and  $\sigma^2$  and can therefore not be used to improve the estimate.

**ANSWER.** Here is the mathematical relationship:

$$(9.3.23) \quad \text{MSE}[\alpha\bar{y}; \mu] = E[(\alpha\bar{y} - \mu)^2] = E[(\alpha\bar{y} - \alpha\mu + \alpha\mu - \mu)^2] < \text{MSE}[\bar{y}; \mu] = \text{var}[\bar{y}]$$

$$(9.3.24) \quad \alpha^2\sigma^2/n + (1 - \alpha)^2\mu^2 < \sigma^2/n$$

Now simplify it:

$$(9.3.25) \quad (1 - \alpha)^2\mu^2 < (1 - \alpha^2)\sigma^2/n = (1 - \alpha)(1 + \alpha)\sigma^2/n$$

This cannot be true for  $\alpha \geq 1$ , because for  $\alpha = 1$  one has equality, and for  $\alpha > 1$ , the righthand side is negative. Therefore we are allowed to assume  $\alpha < 1$ , and can divide by  $1 - \alpha$  without disturbing the inequality:

$$(9.3.26) \quad (1 - \alpha)\mu^2 < (1 + \alpha)\sigma^2/n$$

$$(9.3.27) \quad \mu^2 - \sigma^2/n < \alpha(\mu^2 + \sigma^2/n)$$

The answer is therefore

$$(9.3.28) \quad \frac{n\mu^2 - \sigma^2}{n\mu^2 + \sigma^2} < \alpha < 1.$$

This the range. Note that  $n\mu^2 - \sigma^2 < 0$  may be negative. The best value is in the middle of the range, see Problem 165.



PROBLEM 165. [KS79, example 17.14 on p. 22] *The mathematics in the following problem is easier than it looks. If you can't prove a., assume it and derive b. from it, etc.*

• a. 2 points Let  $\mathbf{t}$  be an estimator of the nonrandom scalar parameter  $\theta$ .  $E[\mathbf{t} - \theta]$  is called the bias of  $\mathbf{t}$ , and  $E[(\mathbf{t} - \theta)^2]$  is called the mean squared error of  $\mathbf{t}$  as an estimator of  $\theta$ , written  $\text{MSE}[\mathbf{t}; \theta]$ . Show that the MSE is the variance plus the squared bias, i.e., that

$$(9.3.29) \quad \text{MSE}[\mathbf{t}; \theta] = \text{var}[\mathbf{t}] + (E[\mathbf{t} - \theta])^2.$$

ANSWER. The most elegant proof, which also indicates what to do when  $\theta$  is random, is:

$$(9.3.30) \quad \text{MSE}[\mathbf{t}; \theta] = E[(\mathbf{t} - \theta)^2] = \text{var}[\mathbf{t} - \theta] + (E[\mathbf{t} - \theta])^2 = \text{var}[\mathbf{t}] + (E[\mathbf{t} - \theta])^2. \quad \square$$

• b. 2 points For the rest of this problem assume that  $\mathbf{t}$  is an unbiased estimator of  $\theta$  with  $\text{var}[\mathbf{t}] > 0$ . We will investigate whether one can get a better MSE if one estimates  $\theta$  by a constant multiple  $a$  instead of  $\mathbf{t}$ . Show that

$$(9.3.31) \quad \text{MSE}[a\mathbf{t}; \theta] = a^2 \text{var}[\mathbf{t}] + (a - 1)^2 \theta^2.$$

ANSWER.  $\text{var}[a\mathbf{t}] = a^2 \text{var}[\mathbf{t}]$  and the bias of  $a\mathbf{t}$  is  $E[a\mathbf{t} - \theta] = (a - 1)\theta$ . Now apply (9.3.30).  $\square$

• c. 1 point Show that, whenever  $a > 1$ , then  $\text{MSE}[a\mathbf{t}; \theta] > \text{MSE}[\mathbf{t}; \theta]$ . If one wants to decrease the MSE, one should therefore not choose  $a > 1$ .

ANSWER.  $\text{MSE}[a\mathbf{t}; \theta] - \text{MSE}[\mathbf{t}; \theta] = (a^2 - 1) \text{var}[\mathbf{t}] + (a - 1)^2 \theta^2 > 0$  since  $a > 1$  and  $\text{var}[\mathbf{t}] > 0$ .  $\square$

• d. 2 points Show that

$$(9.3.32) \quad \left. \frac{d}{da} \text{MSE}[a\mathbf{t}; \theta] \right|_{a=1} > 0.$$

From this follows that the MSE of  $a\mathbf{t}$  is smaller than the MSE of  $\mathbf{t}$ , as long as  $a < 1$  and close enough to 1.

ANSWER. The derivative of (9.3.31) is

$$(9.3.33) \quad \frac{d}{da} \text{MSE}[a\mathbf{t}; \theta] = 2a \text{var}[\mathbf{t}] + 2(a - 1)\theta^2$$

Plug  $a = 1$  into this to get  $2 \text{var}[\mathbf{t}] > 0$ .  $\square$

• e. 2 points By solving the first order condition show that the factor  $a$  which gives smallest MSE is

$$(9.3.34) \quad a = \frac{\theta^2}{\text{var}[\mathbf{t}] + \theta^2}.$$

ANSWER. Rewrite (9.3.33) as  $2a(\text{var}[\mathbf{t}] + \theta^2) - 2\theta^2$  and set it zero.  $\square$

• f. 1 point Assume  $\mathbf{t}$  has an exponential distribution with parameter  $\lambda > 0$ , i.e.,

$$(9.3.35) \quad f_{\mathbf{t}}(t) = \lambda \exp(-\lambda t), \quad t \geq 0 \quad \text{and} \quad f_{\mathbf{t}}(t) = 0 \quad \text{otherwise.}$$

Check that  $f_{\mathbf{t}}(t)$  is indeed a density function.

ANSWER. Since  $\lambda > 0$ ,  $f_{\mathbf{t}}(t) > 0$  for all  $t \geq 0$ . To evaluate  $\int_0^\infty \lambda \exp(-\lambda t) dt$ , substitute  $s = -\lambda t$ , therefore  $ds = -\lambda dt$ , and the upper integration limit changes from  $+\infty$  to  $-\infty$ , therefore the integral is  $-\int_0^{-\infty} \exp(s) ds = 1$ .

• g. 4 points Using this density function (and no other knowledge about exponential distribution) prove that  $\mathbf{t}$  is an unbiased estimator of  $1/\lambda$ , with  $\text{var}[\mathbf{t}] = 1/\lambda^2$ .

ANSWER. To evaluate  $\int_0^\infty \lambda t \exp(-\lambda t) dt$ , use partial integration  $\int uv' dt = uv - \int u'v dt$  with  $u = t$ ,  $u' = 1$ ,  $v = -\exp(-\lambda t)$ ,  $v' = \lambda \exp(-\lambda t)$ . Therefore the integral is  $-t \exp(-\lambda t) \Big|_0^\infty + \int_0^\infty \exp(-\lambda t) dt = 1/\lambda$ , since we just saw that  $\int_0^\infty \lambda \exp(-\lambda t) dt = 1$ .

To evaluate  $\int_0^\infty \lambda t^2 \exp(-\lambda t) dt$ , use partial integration with  $u = t^2$ ,  $u' = 2t$ ,  $v = -\exp(-\lambda t)$ ,  $v' = \lambda \exp(-\lambda t)$ . Therefore the integral is  $-t^2 \exp(-\lambda t) \Big|_0^\infty + 2 \int_0^\infty t \exp(-\lambda t) dt = \frac{2}{\lambda} \int_0^\infty \lambda t \exp(-\lambda t) dt = 2/\lambda^2$ . Therefore  $\text{var}[\mathbf{t}] = E[\mathbf{t}^2] - (E[\mathbf{t}])^2 = 2/\lambda^2 - 1/\lambda^2 = 1/\lambda^2$ .

• h. 2 points Which multiple of  $\mathbf{t}$  has the lowest MSE as an estimator of  $1/\lambda$ ?

ANSWER. It is  $\mathbf{t}/2$ . Just plug  $\theta = 1/\lambda$  into (9.3.34).

$$(9.3.36) \quad a = \frac{1/\lambda^2}{\text{var}[\mathbf{t}] + 1/\lambda^2} = \frac{1/\lambda^2}{1/\lambda^2 + 1/\lambda^2} = \frac{1}{2}.$$

• i. 2 points Assume  $\mathbf{t}_1, \dots, \mathbf{t}_n$  are independently distributed, and each of them has the exponential distribution with the same parameter  $\lambda$ . Which multiple of the sample mean  $\bar{\mathbf{t}} = \frac{1}{n} \sum_{i=1}^n \mathbf{t}_i$  has best MSE as estimator of  $1/\lambda$ ?

ANSWER.  $\bar{\mathbf{t}}$  has expected value  $1/\lambda$  and variance  $1/n\lambda^2$ . Therefore

$$(9.3.37) \quad a = \frac{1/\lambda^2}{\text{var}[\bar{\mathbf{t}}] + 1/\lambda^2} = \frac{1/\lambda^2}{1/n\lambda^2 + 1/\lambda^2} = \frac{n}{n+1},$$

i.e., for the best estimator  $\bar{\mathbf{t}}$  divide the sum by  $n + 1$  instead of  $n$ .

• j. 3 points Assume  $\mathbf{q} \sim \sigma^2 \chi_m^2$  (in other words,  $\frac{1}{\sigma^2} \mathbf{q} \sim \chi_m^2$ , a Chi-square distribution with  $m$  degrees of freedom). Using the fact that  $E[\chi_m^2] = m$  and  $\text{var}[\chi_m^2] = 2m$  compute that multiple of  $\mathbf{q}$  that has minimum MSE as estimator of  $\sigma^2$ .

ANSWER. This is a trick question since  $\mathbf{q}$  itself is not an unbiased estimator of  $\sigma^2$ .  $E[\mathbf{q}] = m\sigma^2$  therefore  $\mathbf{q}/m$  is the unbiased estimator. Since  $\text{var}[\mathbf{q}/m] = 2\sigma^4/m$ , it follows from (9.3.34) that  $a = m/(m + 2)$ , therefore the minimum MSE multiple of  $\mathbf{q}$  is  $\frac{\mathbf{q}}{m} \frac{m}{m+2} = \frac{\mathbf{q}}{m+2}$ . I.e., divide  $\mathbf{q}$  by  $m + 2$  instead of  $m$ .

• k. 3 points Assume you have  $n$  independent observations of a Normally distributed random variable  $\mathbf{y}$  with unknown mean  $\mu$  and standard deviation  $\sigma^2$ . The best unbiased estimator of  $\sigma^2$  is  $\frac{1}{n-1} \sum (\mathbf{y}_i - \bar{\mathbf{y}})^2$ , and the maximum likelihood estimator is  $\frac{1}{n} \sum (\mathbf{y}_i - \bar{\mathbf{y}})^2$ . What are the implications of the above for the question whether one should use the first or the second or still some other multiple of  $\sum (\mathbf{y}_i - \bar{\mathbf{y}})^2$ ?

ANSWER. Taking that multiple of the sum of squared errors which makes the estimator unbiased is not necessarily a good choice. In terms of MSE, the best multiple of  $\sum (\mathbf{y}_i - \bar{\mathbf{y}})^2$  is  $\frac{1}{n+1} \sum (\mathbf{y}_i - \bar{\mathbf{y}})^2$ .  $\square$

• l. 3 points We are still in the model defined in k. Which multiple of the sample mean  $\bar{\mathbf{y}}$  has smallest MSE as estimator of  $\mu$ ? How does this example differ from the ones given above? Can this formula have practical significance?

ANSWER. Here the optimal  $a = \frac{\mu^2}{\mu^2 + (\sigma^2/n)}$ . Unlike in the earlier examples, this  $a$  depends on the unknown parameters. One can “operationalize” it by estimating the parameters from the data, but the noise introduced by this estimation can easily make the estimator worse than the simple  $\bar{\mathbf{y}}$ . Indeed,  $\bar{\mathbf{y}}$  is admissible, i.e., it cannot be uniformly improved upon. On the other hand, the Stein rule, which can be considered an operationalization of a very similar formula (the only difference being that one estimates the mean vector of a vector with at least 3 elements), by estimating  $\mu^2$  and  $\mu^2 + \frac{1}{n}\sigma^2$  from the data, shows that such an operationalization is sometimes successful.  $\square$

We will discuss here one more property of  $\bar{\mathbf{y}}$  and  $s_u^2$ : They together form sufficient statistics for  $\mu$  and  $\sigma^2$ . I.e., any estimator of  $\mu$  and  $\sigma^2$  which is not a function of  $\bar{\mathbf{y}}$  and  $s_u^2$  is less efficient than it could be. Since the factorization theorem for sufficient statistics holds even if the parameter  $\theta$  and its estimate  $\mathbf{t}$  are vectors, we have to write the joint density of the observation vector  $\mathbf{y}$  as a product of two functions, one depending on the parameters and the sufficient statistics, and the other depending on the value taken by  $\mathbf{y}$ , but not on the parameters. Indeed, it will turn out that this second function can just be taken to be  $h(\mathbf{y}) = 1$ , since the density function can be rearranged as

$$(9.3.38) \quad f_{\mathbf{y}}(y_1, \dots, y_n; \mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp\left(-\sum_{i=1}^n (y_i - \mu)^2 / 2\sigma^2\right) =$$

$$(9.3.39) \quad = (2\pi\sigma^2)^{-n/2} \exp\left(-\left(\sum_{i=1}^n (y_i - \bar{y})^2 - n(\bar{y} - \mu)^2\right) / 2\sigma^2\right) =$$

$$(9.3.40) \quad = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{(n-1)s_u^2 - n(\bar{y} - \mu)^2}{2\sigma^2}\right).$$

## CHAPTER 10

## Estimation Principles and Classification of Estimators

### 10.1. Asymptotic or Large-Sample Properties of Estimators

We will discuss asymptotic properties first, because the idea of estimation is to get more certainty by increasing the sample size.

Strictly speaking, asymptotic properties do not refer to individual estimators but to sequences of estimators, one for each sample size  $n$ . And strictly speaking, if one alters the first 10 estimators or the first million estimators and leaves the others unchanged, one still gets a sequence with the same asymptotic properties. The results that follow should therefore be used with caution. The asymptotic properties may say very little about the concrete estimator at hand.

The most basic asymptotic property is (weak) consistency. An estimator  $\mathbf{t}_n$  (where  $n$  is the sample size) of the parameter  $\theta$  is consistent iff

$$(10.1.1) \quad \text{plim}_{n \rightarrow \infty} \mathbf{t}_n = \theta.$$

Roughly, a consistent estimation procedure is one which gives the correct parameter values if the sample is large enough. There are only very few exceptional situations in which an estimator is acceptable which is not *consistent*, i.e., which does not converge in the plim to the true parameter value.

**PROBLEM 166.** *Can you think of a situation where an estimator which is not consistent is acceptable?*

**ANSWER.** If additional data no longer give information, like when estimating the initial state of a timeseries, or in prediction. And if there is no identification but the value can be confined to an interval. This is also inconsistency.  $\square$

The following is an important property of consistent estimators:

*Slutsky theorem:* If  $\mathbf{t}$  is a consistent estimator for  $\theta$ , and the function  $g$  is continuous at the true value of  $\theta$ , then  $g(\mathbf{t})$  is consistent for  $g(\theta)$ .

For the proof of the Slutsky theorem remember the definition of a continuous function.  $g$  is continuous at  $\theta$  iff for all  $\varepsilon > 0$  there exists a  $\delta > 0$  with the property

that for all  $\theta_1$  with  $|\theta_1 - \theta| < \delta$  follows  $|g(\theta_1) - g(\theta)| < \varepsilon$ . To prove consistency of  $g(\mathbf{t})$  we have to show that for all  $\varepsilon > 0$ ,  $\Pr[|g(\mathbf{t}) - g(\theta)| \geq \varepsilon] \rightarrow 0$ . Choose for  $\mathbf{t}$  given  $\varepsilon$  a  $\delta$  as above, then  $|g(\mathbf{t}) - g(\theta)| \geq \varepsilon$  implies  $|\mathbf{t} - \theta| \geq \delta$ , because all the values of  $\mathbf{t}$  for which  $|\mathbf{t} - \theta| < \delta$  lead to a  $g(\mathbf{t})$  with  $|g(\mathbf{t}) - g(\theta)| < \varepsilon$ . This logical implication means that

$$(10.1.2) \quad \Pr[|g(\mathbf{t}) - g(\theta)| \geq \varepsilon] \leq \Pr[|\mathbf{t} - \theta| \geq \delta].$$

Since the probability on the righthand side converges to zero, the one on the lefthand side converges too.

Different consistent estimators can have quite different speeds of convergence. Are there estimators which have optimal asymptotic properties among all consistent estimators? Yes, if one limits oneself to a fairly reasonable subclass of consistent estimators.

Here are the details: Most consistent estimators we will encounter are asymptotically normal, i.e., the “shape” of their distribution function converges toward the normal distribution, as we had it for the sample mean in the central limit theorem. In order to be able to use this asymptotic distribution for significance tests and confidence intervals, however, one needs more than asymptotic normality (as many textbooks are not aware of this): one needs the convergence to normality to be *uniform in compact intervals* [Rao73, p. 346–351]. Such estimators are called *consistent uniformly asymptotically normal* estimators (CUAN estimators)

If one limits oneself to CUAN estimators it can be shown that there are asymptotically “best” CUAN estimators. Since the distribution is asymptotically normal there is no problem to define what it means to be asymptotically best: those estimators are asymptotically best whose asymptotic MSE = asymptotic variance is smallest. CUAN estimators whose MSE is asymptotically no larger than that of any other CUAN estimator, are called *asymptotically efficient*. Rao has shown that for CUAN estimators the lower bound for this asymptotic variance is the asymptotic limit of the Cramer Rao lower bound (CRLB). (More about the CRLB below). Maximum likelihood estimators are therefore usually efficient CUAN estimators. In this sense one can think of maximum likelihood estimators to be something like asymptotically best consistent estimators, compare a statement to this effect in [Ame94, 144]. And one can think of asymptotically efficient CUAN estimators as estimators who are in large samples as good as maximum likelihood estimators.

All these are large sample properties. Among the asymptotically efficient estimators there are still wide differences regarding the small sample properties. Asymptotic efficiency should therefore again be considered a minimum requirement: there may be very good reasons *not* to be working with an asymptotically efficient estimator.

**PROBLEM 167.** *Can you think of situations in which an estimator is acceptable which is not asymptotically efficient?*

ANSWER. If robustness matters then the median may be preferable to the mean, although it is less efficient.  $\square$

## 10.2. Small Sample Properties

In order to judge how good an estimator is for small samples, one has two dilemmas: (1) there are many different criteria for an estimator to be “good”; (2) even if one has decided on one criterion, a given estimator may be good for some values of the unknown parameters and not so good for others.

If  $\mathbf{x}$  and  $\mathbf{y}$  are two estimators of the parameter  $\theta$ , then each of the following conditions can be interpreted to mean that  $\mathbf{x}$  is better than  $\mathbf{y}$ :

$$(10.2.1) \quad \Pr[|\mathbf{x} - \theta| \leq |\mathbf{y} - \theta|] = 1$$

$$(10.2.2) \quad E[g(\mathbf{x} - \theta)] \leq E[g(\mathbf{y} - \theta)]$$

for every continuous function  $g$  which is and nonincreasing for  $x < 0$  and nondecreasing for  $x > 0$

$$(10.2.3) \quad E[g(|\mathbf{x} - \theta|)] \leq E[g(|\mathbf{y} - \theta|)]$$

for every continuous and nondecreasing function  $g$

$$(10.2.4) \quad \Pr\{|\mathbf{x} - \theta| > \varepsilon\} \leq \Pr\{|\mathbf{y} - \theta| > \varepsilon\} \quad \text{for every } \varepsilon$$

$$(10.2.5) \quad E[(\mathbf{x} - \theta)^2] \leq E[(\mathbf{y} - \theta)^2]$$

$$(10.2.6) \quad \Pr[|\mathbf{x} - \theta| < |\mathbf{y} - \theta|] \geq \Pr[|\mathbf{x} - \theta| > |\mathbf{y} - \theta|]$$

This list is from [Ame94, pp. 118–122]. But we will simply use the MSE.

Therefore we are left with dilemma (2). There is no single estimator that has uniformly the smallest MSE in the sense that its MSE is better than the MSE of *any* other estimator whatever the value of the parameter value. To see this, simply think of the following estimator  $\mathbf{t}$  of  $\theta$ :  $\mathbf{t} = 10$ ; i.e., whatever the outcome of the experiments,  $\mathbf{t}$  always takes the value 10. This estimator has zero MSE when  $\theta$  happens to be 10, but is a bad estimator when  $\theta$  is far away from 10. If an estimator existed which had uniformly best MSE, then it had to be better than all the constant estimators, i.e., have zero MSE whatever the value of the parameter, and this is only possible if the parameter itself is observed.

Although the MSE criterion cannot be used to pick one best estimator, it can be used to rule out estimators which are unnecessarily bad in the sense that other estimators exist which are never worse but sometimes better in terms of MSE whatever the true parameter values. Estimators which are dominated in this sense are called inadmissible.

But how can one choose between two admissible estimators? [Ame94, p. 124] gives two reasonable strategies. One is to integrate the MSE out over a distribution

of the likely values of the parameter. This is in the spirit of the Bayesians, although Bayesians would still do it differently. The other strategy is to choose a minimum variance unbiased strategy. Amemiya seems to consider this an alright strategy, but it is really too defensive. Here is a third strategy, which is often used but less well founded theoretically: Since there are no estimators which have minimum MSE among all estimators, one often looks for estimators which have minimum MSE among all estimators with a certain property. And the “certain property” which is most often used is unbiasedness. The MSE of an unbiased estimator is its variance; and an estimator which has minimum variance in the class of all unbiased estimators is called “efficient.”

The class of unbiased estimators has a high-sounding name, and the results related with Cramer-Rao and Least Squares seem to confirm that it is an important class of estimators. However I will argue in these class notes that unbiasedness itself is not a desirable property.

## 10.3. Comparison Unbiasedness Consistency

Let us compare consistency with unbiasedness. If the estimator is unbiased then its expected value for any sample size, whether large or small, is equal to the true parameter value. By the law of large numbers this can be translated into a statement about large samples: The mean of many independent replications of the estimate, *even if each replication only uses a small number of observations*, gives the true parameter value. Unbiasedness says therefore something about the small sample properties of the estimator, while consistency does not.

The following thought experiment may clarify the difference between unbiasedness and consistency. Imagine you are conducting an experiment which gives you every ten seconds an independent measurement, i.e., a measurement whose value is not influenced by the outcome of previous measurements. Imagine further that the experimental setup is connected to a computer which estimates certain parameters from that experiment, re-calculating its estimate every time twenty new observations have become available, and which displays the current values of the estimate on a screen. And assume that the estimation procedure used by the computer is consistent, but biased for any finite number of observations.

Consistency means: after a sufficiently long time, the digits of the parameter estimate displayed by the computer will be correct. That the estimator is biased means: if the computer were to use every batch of 20 observations to form a new estimate of the parameter, without utilizing prior observations, and then would update the *average* of all these independent estimates as its updated estimate, it would end up displaying a wrong parameter value on the screen.

A biased estimator gives, even in the limit, an incorrect result as long as one uses an updating procedure is the simple taking the averages of all previous estimates. If an estimator is biased but consistent, then a better updating method is available.

which will end up in the correct parameter value. A biased estimator therefore is not necessarily one which gives incorrect information about the parameter value; but it is one which one cannot update by simply taking averages. But there is no reason to limit oneself to such a crude method of updating. Obviously the question whether the estimate is biased is of little relevance, as long as it is consistent. The moral of the story is: If one looks for desirable estimators, by no means should one restrict one's search to unbiased estimators! The high-sounding name "unbiased" for the technical property  $E[\mathbf{t}] = \theta$  has created a lot of confusion.

Besides having no advantages, the category of unbiasedness even has some inconvenient properties: In some cases, in which consistent estimators exist, there are no unbiased estimators. And if an estimator  $\mathbf{t}$  is an unbiased estimate for the parameter  $\theta$ , then the estimator  $g(\mathbf{t})$  is usually no longer an unbiased estimator for  $g(\theta)$ . It depends on the way a certain quantity is measured whether the estimator is unbiased or not. However consistency carries over.

Unbiasedness is not the only possible criterion which ensures that the values of the estimator are centered over the value it estimates. Here is another plausible definition:

DEFINITION 10.3.1. An estimator  $\hat{\theta}$  of the scalar  $\theta$  is called *median unbiased* for all  $\theta \in \Theta$  iff

$$(10.3.1) \quad \Pr[\hat{\theta} < \theta] = \Pr[\hat{\theta} > \theta] = \frac{1}{2}$$

This concept is always applicable, even for estimators whose expected value does not exist.

PROBLEM 168. *6 points (Not eligible for in-class exams) The purpose of the following problem is to show how restrictive the requirement of unbiasedness is. Sometimes no unbiased estimators exist, and sometimes, as in the example here, unbiasedness leads to absurd estimators. Assume the random variable  $\mathbf{x}$  has the geometric distribution with parameter  $p$ , where  $0 \leq p \leq 1$ . In other words, it can only assume the integer values  $1, 2, 3, \dots$ , with probabilities*

$$(10.3.2) \quad \Pr[\mathbf{x} = r] = (1 - p)^{r-1}p.$$

Show that the unique unbiased estimator of  $p$  on the basis of one observation of  $\mathbf{x}$  is the random variable  $f(\mathbf{x})$  defined by  $f(x) = 1$  if  $x = 1$  and  $0$  otherwise. Hint: Use the mathematical fact that a function  $\phi(q)$  that can be expressed as a power series  $\phi(q) = \sum_{j=0}^{\infty} a_j q^j$ , and which takes the values  $\phi(q) = 1$  for all  $q$  in some interval of nonzero length, is the power series with  $a_0 = 1$  and  $a_j = 0$  for  $j \neq 0$ . (You will need the hint at the end of your answer, don't try to start with the hint!)

ANSWER. Unbiasedness means that  $E[f(\mathbf{x})] = \sum_{r=1}^{\infty} f(r)(1-p)^{r-1}p = p$  for all  $p$  in the unit interval, therefore  $\sum_{r=1}^{\infty} f(r)(1-p)^{r-1} = 1$ . This is a power series in  $q = 1-p$ , which must be

identically equal to 1 for all values of  $q$  between 0 and 1. An application of the hint shows that the constant term in this power series, corresponding to the value  $r-1=0$ , must be  $=1$ , and other  $f(r) = 0$ . Here older formulation: An application of the hint with  $q = 1-p$ ,  $j = r-1$ ,  $a_j = f(j+1)$  gives  $f(1) = 1$  and all other  $f(r) = 0$ . This estimator is absurd since it lies on the boundary of the range of possible values for  $q$ .

PROBLEM 169. *As in Question 61, you make two independent trials of a Bernoulli experiment with success probability  $\theta$ , and you observe  $\mathbf{t}$ , the number of successes*

- a. Give an unbiased estimator of  $\theta$  based on  $\mathbf{t}$  (i.e., which is a function of  $\mathbf{t}$ ).
- b. Give an unbiased estimator of  $\theta^2$ .
- c. Show that there is no unbiased estimator of  $\theta^3$ .

Hint: Since  $\mathbf{t}$  can only take the three values 0, 1, and 2, any estimator  $\mathbf{u}$  which is a function of  $\mathbf{t}$  is determined by the values it takes when  $\mathbf{t}$  is 0, 1, or 2, call them  $u_0, u_1$ , and  $u_2$ . Express  $E[\mathbf{u}]$  as a function of  $u_0, u_1$ , and  $u_2$ .

ANSWER.  $E[\mathbf{u}] = u_0(1-\theta)^2 + 2u_1\theta(1-\theta) + u_2\theta^2 = u_0 + (2u_1 - 2u_0)\theta + (u_0 - 2u_1 + u_2)\theta^2$ . This is always a second degree polynomial in  $\theta$ , therefore whatever is not a second degree polynomial cannot be the expected value of any function of  $\mathbf{t}$ . For  $E[\mathbf{u}] = \theta$  we need  $u_0 = 0$ ,  $2u_1 - 2u_0 = 2u_1 = 1$ , therefore  $u_1 = 0.5$ , and  $u_0 - 2u_1 + u_2 = -1 + u_2 = 0$ , i.e.  $u_2 = 1$ . This is, in other words,  $\mathbf{u} = \mathbf{t}$ . For  $E[\mathbf{u}] = \theta^2$  we need  $u_0 = 0$ ,  $2u_1 - 2u_0 = 2u_1 = 0$ , therefore  $u_1 = 0$ , and  $u_0 - 2u_1 + u_2 = u_2 = 1$ . This is, in other words,  $\mathbf{u} = \mathbf{t}(\mathbf{t}-1)/2$ . From this equation one also sees that  $\theta^3$  and higher powers or things like  $1/\theta$ , cannot be the expected values of any estimators.

- d. Compute the moment generating function of  $\mathbf{t}$ .

ANSWER.

$$(10.3.3) \quad E[e^{\lambda \mathbf{t}}] = e^0 \cdot (1-\theta)^2 + e^{\lambda} \cdot 2\theta(1-\theta) + e^{2\lambda} \cdot \theta^2 = (1-\theta + \theta e^{\lambda})^2$$

PROBLEM 170. *This is [KS79, Question 17.11 on p. 34], originally [Fis, p. 70]*

• a. *1 point Assume  $\mathbf{t}$  and  $\mathbf{u}$  are two unbiased estimators of the same unknown scalar nonrandom parameter  $\theta$ .  $\mathbf{t}$  and  $\mathbf{u}$  have finite variances and satisfy  $\text{var}[\mathbf{u} - \mathbf{t}] > 0$ . Show that a linear combination of  $\mathbf{t}$  and  $\mathbf{u}$ , i.e., an estimator of  $\theta$  which can be written in the form  $\alpha \mathbf{t} + \beta \mathbf{u}$ , is unbiased if and only if  $\alpha = 1 - \beta$ . In other words, any unbiased estimator which is a linear combination of  $\mathbf{t}$  and  $\mathbf{u}$  can be written in the form*

$$(10.3.4) \quad \mathbf{t} + \beta(\mathbf{u} - \mathbf{t}).$$

• b. *2 points By solving the first order condition show that the unbiased linear combination of  $\mathbf{t}$  and  $\mathbf{u}$  which has lowest MSE is*

$$(10.3.5) \quad \hat{\theta} = \mathbf{t} - \frac{\text{cov}[\mathbf{t}, \mathbf{u} - \mathbf{t}]}{\text{var}[\mathbf{u} - \mathbf{t}]}(\mathbf{u} - \mathbf{t})$$

*Hint: your arithmetic will be simplest if you start with (10.3.4).*

- c. 1 point If  $\rho^2$  is the squared correlation coefficient between  $\mathbf{t}$  and  $\mathbf{u} - \mathbf{t}$ , i.e.,

$$(10.3.6) \quad \rho^2 = \frac{(\text{cov}[\mathbf{t}, \mathbf{u} - \mathbf{t}])^2}{\text{var}[\mathbf{t}] \text{var}[\mathbf{u} - \mathbf{t}]}$$

show that  $\text{var}[\hat{\theta}] = \text{var}[\mathbf{t}](1 - \rho^2)$ .

- d. 1 point Show that  $\text{cov}[\mathbf{t}, \mathbf{u} - \mathbf{t}] \neq 0$  implies  $\text{var}[\mathbf{u} - \mathbf{t}] \neq 0$ .
- e. 2 points Use (10.3.5) to show that if  $\mathbf{t}$  is the minimum MSE unbiased estimator of  $\theta$ , and  $\mathbf{u}$  another unbiased estimator of  $\theta$ , then

$$(10.3.7) \quad \text{cov}[\mathbf{t}, \mathbf{u} - \mathbf{t}] = 0.$$

- f. 1 point Use (10.3.5) to show also the opposite: if  $\mathbf{t}$  is an unbiased estimator of  $\theta$  with the property that  $\text{cov}[\mathbf{t}, \mathbf{u} - \mathbf{t}] = 0$  for every other unbiased estimator  $\mathbf{u}$  of  $\theta$ , then  $\mathbf{t}$  has minimum MSE among all unbiased estimators of  $\theta$ .

There are estimators which are consistent but their bias does not converge to zero:

$$(10.3.8) \quad \hat{\theta}_n = \begin{cases} \theta & \text{with probability } 1 - \frac{1}{n} \\ n & \text{with probability } \frac{1}{n} \end{cases}$$

Then  $\Pr(|\hat{\theta}_n - \theta| \geq \varepsilon) \leq \frac{1}{n}$ , i.e., the estimator is consistent, but  $E[\hat{\theta}] = \theta \frac{n-1}{n} + n \rightarrow \theta + 1 \neq \theta$ .

**PROBLEM 171.** 4 points Is it possible to have a consistent estimator whose bias becomes unbounded as the sample size increases? Either prove that it is not possible or give an example.

**ANSWER.** Yes, this can be achieved by making the rare outliers even wilder than in (10.3.8), say

$$(10.3.9) \quad \hat{\theta}_n = \begin{cases} \theta & \text{with probability } 1 - \frac{1}{n} \\ n^2 & \text{with probability } \frac{1}{n} \end{cases}$$

Here  $\Pr(|\hat{\theta}_n - \theta| \geq \varepsilon) \leq \frac{1}{n}$ , i.e., the estimator is consistent, but  $E[\hat{\theta}] = \theta \frac{n-1}{n} + n \rightarrow \theta + n$ .  $\square$

And of course there are estimators which are unbiased but not consistent: simply take the first observation  $x_1$  as an estimator if  $E[x]$  and ignore all the other observations.

## 10.4. The Cramer-Rao Lower Bound

Take a scalar random variable  $\mathbf{y}$  with density function  $f_{\mathbf{y}}$ . The entropy of  $\mathbf{y}$ , if it exists, is  $H[\mathbf{y}] = -E[\log(f_{\mathbf{y}}(\mathbf{y}))]$ . This is the continuous equivalent of (3.11.2). The entropy is the measure of the amount of randomness in this variable. If there is little information and much noise in this variable, the entropy is high.

Now let  $y \mapsto g(y)$  be the density function of a different random variable  $\mathbf{x}$ . In other words,  $g$  is some function which satisfies  $g(y) \geq 0$  for all  $y$ , and  $\int_{-\infty}^{+\infty} g(y) dy = 1$ . Equation (3.11.10) with  $v = g(y)$  and  $w = f_{\mathbf{y}}(y)$  gives

$$(10.4.1) \quad f_{\mathbf{y}}(y) - f_{\mathbf{y}}(y) \log f_{\mathbf{y}}(y) \leq g(y) - f_{\mathbf{y}}(y) \log g(y).$$

This holds for every value  $y$ , and integrating over  $y$  gives  $1 - E[\log f_{\mathbf{y}}(\mathbf{y})] \leq 1 - E[\log g(\mathbf{y})]$  or

$$(10.4.2) \quad E[\log f_{\mathbf{y}}(\mathbf{y})] \geq E[\log g(\mathbf{y})].$$

This is an important extremal value property which distinguishes the density function  $f_{\mathbf{y}}(y)$  of  $\mathbf{y}$  from all other density functions: That density function  $g$  which maximizes  $E[\log g(\mathbf{y})]$  is  $g = f_{\mathbf{y}}$ , the true density function of  $\mathbf{y}$ .

This optimality property lies at the basis of the Cramer-Rao inequality, and is also the reason why maximum likelihood estimation is so good. The difference between the left and right hand side in (10.4.2) is called the Kullback-Leibler discrepancy between the random variables  $\mathbf{y}$  and  $\mathbf{x}$  (where  $\mathbf{x}$  is a random variable whose density is  $g$ ).

The Cramer Rao inequality gives a lower bound for the MSE of an unbiased estimator of the parameter of a probability distribution (which has to satisfy certain regularity conditions). This allows one to determine whether a given unbiased estimator has a MSE as low as any other unbiased estimator (i.e., whether it is "efficient.")

**PROBLEM 172.** Assume the density function of  $\mathbf{y}$  depends on a parameter  $\theta$ . Write it  $f_{\mathbf{y}}(y; \theta)$ , and  $\theta^\circ$  is the true value of  $\theta$ . In this problem we will compare the expected value of  $\mathbf{y}$  and of functions of  $\mathbf{y}$  with what would be their expected values if the true parameter value were not  $\theta^\circ$  but would take some other value  $\theta$ . If  $\mathbf{t}$  is a random variable which is a function of  $\mathbf{y}$ , we write  $E_\theta[\mathbf{t}]$  for what would be the expected value of  $\mathbf{t}$  if the true value of the parameter were  $\theta$  instead of  $\theta^\circ$ . Occasionally, we will use the subscript  $\circ$  as in  $E_\circ$  to indicate that we are dealing here with the usual case in which the expected value is taken with respect to the true parameter value  $\theta^\circ$ . Instead of  $E_\circ$  one usually simply writes  $E$ , since it is usually self-understood that one has to plug the right parameter values into the density function if one takes expected values. The subscript  $\circ$  is necessary here only because in the present problem,

sometimes take expected values with respect to the “wrong” parameter values. The same notational convention also applies to variances, covariances, and the MSE.

Throughout this problem we assume that the following regularity conditions hold: (a) the range of  $\mathbf{y}$  is independent of  $\theta$ , and (b) the derivative of the density function with respect to  $\theta$  is a continuous differentiable function of  $\theta$ . These regularity conditions ensure that one can differentiate under the integral sign, i.e., for all function  $t(y)$  follows

$$(10.4.3) \quad \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_{\mathbf{y}}(y; \theta) t(y) dy = \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f_{\mathbf{y}}(y; \theta) t(y) dy = \frac{\partial}{\partial \theta} E_{\theta}[t(\mathbf{y})]$$

$$(10.4.4) \quad \int_{-\infty}^{\infty} \frac{\partial^2}{(\partial \theta)^2} f_{\mathbf{y}}(y; \theta) t(y) dy = \frac{\partial^2}{(\partial \theta)^2} \int_{-\infty}^{\infty} f_{\mathbf{y}}(y; \theta) t(y) dy = \frac{\partial^2}{(\partial \theta)^2} E_{\theta}[t(\mathbf{y})].$$

- a. 1 point The score is defined as the random variable

$$(10.4.5) \quad q(\mathbf{y}; \theta) = \frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta).$$

In other words, we do three things to the density function: take its logarithm, then take the derivative of this logarithm with respect to the parameter, and then plug the random variable into it. This gives us a random variable which also depends on the nonrandom parameter  $\theta$ . Show that the score can also be written as

$$(10.4.6) \quad q(\mathbf{y}; \theta) = \frac{1}{f_{\mathbf{y}}(\mathbf{y}; \theta)} \frac{\partial f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta}$$

ANSWER. This is the chain rule for differentiation: for any differentiable function  $g(\theta)$ ,  $\frac{\partial}{\partial \theta} \log g(\theta) = \frac{1}{g(\theta)} \frac{\partial g(\theta)}{\partial \theta}$ . □

- b. 1 point If the density function is member of an exponential dispersion family (??), show that the score function has the form

$$(10.4.7) \quad q(\mathbf{y}; \theta) = \frac{\mathbf{y} - \frac{\partial b(\theta)}{\partial \theta}}{a(\psi)}$$

ANSWER. This is a simple substitution: if

$$(10.4.8) \quad f_{\mathbf{y}}(y; \theta, \psi) = \exp\left(\frac{y\theta - b(\theta)}{a(\psi)} + c(y, \psi)\right),$$

then

$$(10.4.9) \quad \frac{\partial \log f_{\mathbf{y}}(\mathbf{y}; \theta, \psi)}{\partial \theta} = \frac{\mathbf{y} - \frac{\partial b(\theta)}{\partial \theta}}{a(\psi)}$$

□

- c. 3 points If  $f_{\mathbf{y}}(y; \theta^\circ)$  is the true density function of  $\mathbf{y}$ , then we know from (10.4.2) that  $E_{\circ}[\log f_{\mathbf{y}}(\mathbf{y}; \theta^\circ)] \geq E_{\circ}[\log f(\mathbf{y}; \theta)]$  for all  $\theta$ . This explains why the score is so important: it is the derivative of that function whose expected value is maximized if the true parameter is plugged into the density function. The first-order conditions in this situation read: the expected value of this derivative must be zero for the true parameter value. This is the next thing you are asked to show: If  $\theta^\circ$  is the true parameter value, show that  $E_{\circ}[q(\mathbf{y}; \theta^\circ)] = 0$ .

ANSWER. First write for general  $\theta$

$$(10.4.10) \quad E_{\circ}[q(\mathbf{y}; \theta)] = \int_{-\infty}^{\infty} q(y; \theta) f_{\mathbf{y}}(y; \theta^\circ) dy = \int_{-\infty}^{\infty} \frac{1}{f_{\mathbf{y}}(y; \theta)} \frac{\partial f_{\mathbf{y}}(y; \theta)}{\partial \theta} f_{\mathbf{y}}(y; \theta^\circ) dy.$$

For  $\theta = \theta^\circ$  this simplifies:

$$(10.4.11) \quad E_{\circ}[q(\mathbf{y}; \theta^\circ)] = \int_{-\infty}^{\infty} \frac{\partial f_{\mathbf{y}}(y; \theta)}{\partial \theta} \Big|_{\theta=\theta^\circ} dy = \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f_{\mathbf{y}}(y; \theta) dy \Big|_{\theta=\theta^\circ} = \frac{\partial}{\partial \theta} 1 = 0.$$

Here I am writing  $\frac{\partial f_{\mathbf{y}}(y; \theta)}{\partial \theta} \Big|_{\theta=\theta^\circ}$  instead of the simpler notation  $\frac{\partial f_{\mathbf{y}}(y; \theta^\circ)}{\partial \theta}$ , in order to emphasize that one first has to take a derivative with respect to  $\theta$  and then one plugs  $\theta^\circ$  into that derivative.

- d. Show that, in the case of the exponential dispersion family,

$$(10.4.12) \quad E_{\circ}[\mathbf{y}] = \frac{\partial b(\theta)}{\partial \theta} \Big|_{\theta=\theta^\circ}$$

ANSWER. Follows from the fact that the score function of the exponential family (10.4.7) has zero expected value.

- e. 5 points If we differentiate the score, we obtain the Hessian

$$(10.4.13) \quad \mathbf{h}(\theta) = \frac{\partial^2}{(\partial \theta)^2} \log f_{\mathbf{y}}(\mathbf{y}; \theta).$$

From now on we will write the score function as  $\mathbf{q}(\theta)$  instead of  $q(\mathbf{y}; \theta)$ ; i.e., we no longer make it explicit that  $\mathbf{q}$  is a function of  $\mathbf{y}$  but write it as a random variable which depends on the parameter  $\theta$ . We also suppress the dependence of  $\mathbf{h}$  on  $\mathbf{y}$ ; our notation  $\mathbf{h}(\theta)$  is short for  $\mathbf{h}(\mathbf{y}; \theta)$ . Since there is only one parameter in the density function, score and Hessian are scalars; but in the general case, the score is a vector and the Hessian a matrix. Show that, for the true parameter value  $\theta^\circ$ , the negative of the expected value of the Hessian equals the variance of the score, i.e., the expected value of the square of the score:

$$(10.4.14) \quad E_{\circ}[\mathbf{h}(\theta^\circ)] = -E_{\circ}[\mathbf{q}^2(\theta^\circ)].$$

ANSWER. Start with the definition of the score

$$(10.4.15) \quad q(\mathbf{y}; \theta) = \frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta) = \frac{1}{f_{\mathbf{y}}(\mathbf{y}; \theta)} \frac{\partial}{\partial \theta} f_{\mathbf{y}}(\mathbf{y}; \theta),$$

and differentiate the rightmost expression one more time:

$$(10.4.16) \quad h(\mathbf{y}; \theta) = \frac{\partial}{\partial \theta} q(\mathbf{y}; \theta) = -\frac{1}{f_{\mathbf{y}}^2(\mathbf{y}; \theta)} \left( \frac{\partial}{\partial \theta} f_{\mathbf{y}}(\mathbf{y}; \theta) \right)^2 + \frac{1}{f_{\mathbf{y}}(\mathbf{y}; \theta)} \frac{\partial^2}{\partial \theta^2} f_{\mathbf{y}}(\mathbf{y}; \theta)$$

$$(10.4.17) \quad = -q^2(\mathbf{y}; \theta) + \frac{1}{f_{\mathbf{y}}(\mathbf{y}; \theta)} \frac{\partial^2}{\partial \theta^2} f_{\mathbf{y}}(\mathbf{y}; \theta)$$

Taking expectations we get

$$(10.4.18) \quad E_{\circ}[h(\mathbf{y}; \theta)] = -E_{\circ}[q^2(\mathbf{y}; \theta)] + \int_{-\infty}^{+\infty} \frac{1}{f_{\mathbf{y}}(\mathbf{y}; \theta)} \left( \frac{\partial^2}{\partial \theta^2} f_{\mathbf{y}}(\mathbf{y}; \theta) \right) f_{\mathbf{y}}(\mathbf{y}; \theta^{\circ}) d\mathbf{y}$$

Again, for  $\theta = \theta^{\circ}$ , we can simplify the integrand and differentiate under the integral sign:

$$(10.4.19) \quad \int_{-\infty}^{+\infty} \frac{\partial^2}{\partial \theta^2} f_{\mathbf{y}}(\mathbf{y}; \theta) d\mathbf{y} = \frac{\partial^2}{\partial \theta^2} \int_{-\infty}^{+\infty} f_{\mathbf{y}}(\mathbf{y}; \theta) d\mathbf{y} = \frac{\partial^2}{\partial \theta^2} 1 = 0.$$

□

- f. Derive from (10.4.14) that, for the exponential dispersion family (??),

$$(10.4.20) \quad \text{var}_{\circ}[\mathbf{y}] = \frac{\partial^2 b(\theta)}{\partial \theta^2} a(\phi) \Big|_{\theta=\theta^{\circ}}$$

ANSWER. Differentiation of (10.4.7) gives  $h(\theta) = -\frac{\partial^2 b(\theta)}{\partial \theta^2} \frac{1}{a(\phi)}$ . This is constant and therefore equal to its own expected value. (10.4.14) says therefore

$$(10.4.21) \quad \frac{\partial^2 b(\theta)}{\partial \theta^2} \Big|_{\theta=\theta^{\circ}} \frac{1}{a(\phi)} = E_{\circ}[q^2(\theta^{\circ})] = \frac{1}{(a(\phi))^2} \text{var}_{\circ}[\mathbf{y}]$$

from which (10.4.20) follows. □

PROBLEM 173.

- a. Use the results from question 172 to derive the following strange and interesting result: for any random variable  $\mathbf{t}$  which is a function of  $\mathbf{y}$ , i.e.,  $\mathbf{t} = t(\mathbf{y})$ , follows  $\text{cov}_{\circ}[\mathbf{q}(\theta^{\circ}), \mathbf{t}] = \frac{\partial}{\partial \theta} E_{\theta}[\mathbf{t}] \Big|_{\theta=\theta^{\circ}}$ .

ANSWER. The following equation holds for all  $\theta$ :

$$(10.4.22) \quad E_{\circ}[\mathbf{q}(\theta)\mathbf{t}] = \int_{-\infty}^{\infty} \frac{1}{f_{\mathbf{y}}(\mathbf{y}; \theta)} \frac{\partial f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta} t(\mathbf{y}) f_{\mathbf{y}}(\mathbf{y}; \theta^{\circ}) d\mathbf{y}$$

If the  $\theta$  in  $\mathbf{q}(\theta)$  is the right parameter value  $\theta^{\circ}$  one can simplify:

$$(10.4.23) \quad E_{\circ}[\mathbf{q}(\theta^{\circ})\mathbf{t}] = \int_{-\infty}^{\infty} \frac{\partial f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta} \Big|_{\theta=\theta^{\circ}} t(\mathbf{y}) d\mathbf{y}$$

$$(10.4.24) \quad = \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f_{\mathbf{y}}(\mathbf{y}; \theta) t(\mathbf{y}) d\mathbf{y} \Big|_{\theta=\theta^{\circ}}$$

$$(10.4.25) \quad = \frac{\partial}{\partial \theta} E_{\theta}[\mathbf{t}] \Big|_{\theta=\theta^{\circ}}$$

This is at the same time the covariance:  $\text{cov}_{\circ}[\mathbf{q}(\theta^{\circ}), \mathbf{t}] = E_{\circ}[\mathbf{q}(\theta^{\circ})\mathbf{t}] - E_{\circ}[\mathbf{q}(\theta^{\circ})] E_{\circ}[\mathbf{t}] = E_{\circ}[\mathbf{q}(\theta^{\circ})\mathbf{t}]$  since  $E_{\circ}[\mathbf{q}(\theta^{\circ})] = 0$ .

*Explanation, nothing to prove here: Now if  $\mathbf{t}$  is an unbiased estimator of whatever the value of  $\theta$ , then it follows  $\text{cov}_{\circ}[\mathbf{q}(\theta^{\circ}), \mathbf{t}] = \frac{\partial}{\partial \theta} \theta = 1$ . From this follows by Cauchy-Schwartz  $\text{var}_{\circ}[\mathbf{t}] \text{var}_{\circ}[\mathbf{q}(\theta^{\circ})] \geq 1$ , or  $\text{var}_{\circ}[\mathbf{t}] \geq 1/\text{var}_{\circ}[\mathbf{q}(\theta^{\circ})]$ . Since  $E_{\circ}[\mathbf{q}(\theta^{\circ})] = 0$ , we know  $\text{var}_{\circ}[\mathbf{q}(\theta^{\circ})] = E_{\circ}[q^2(\theta^{\circ})]$ , and since  $\mathbf{t}$  is unbiased, we know  $\text{var}_{\circ}[\mathbf{t}] = \text{MSE}_{\circ}[\mathbf{t}; \theta^{\circ}]$ . Therefore the Cauchy-Schwartz inequality reads*

$$(10.4.26) \quad \text{MSE}_{\circ}[\mathbf{t}; \theta^{\circ}] \geq 1/E_{\circ}[q^2(\theta^{\circ})].$$

*This is the Cramer-Rao inequality. The inverse of the variance of  $\mathbf{q}(\theta^{\circ})$ ,  $1/\text{var}_{\circ}[\mathbf{q}(\theta^{\circ})] = 1/E_{\circ}[q^2(\theta^{\circ})]$ , is called the Fisher information, written  $I(\theta^{\circ})$ . It is a lower bound on the MSE of any unbiased estimator of  $\theta$ . Because of (10.4.14), the Cramer-Rao inequality can also be written in the form*

$$(10.4.27) \quad \text{MSE}[\mathbf{t}; \theta^{\circ}] \geq -1/E_{\circ}[h(\theta^{\circ})].$$

(10.4.26) and (10.4.27) are usually written in the following form: Assume  $\mathbf{y}$  has density function  $f_{\mathbf{y}}(\mathbf{y}; \theta)$  which depends on the unknown parameter  $\theta$ , and and  $\mathbf{t}(\mathbf{y})$  be any unbiased estimator of  $\theta$ . Then

$$(10.4.28) \quad \text{var}[\mathbf{t}] \geq \frac{1}{E[(\frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta))^2]} = \frac{-1}{E[\frac{\partial^2}{\partial \theta^2} \log f_{\mathbf{y}}(\mathbf{y}; \theta)]}.$$

(Sometimes the first and sometimes the second expression is easier to evaluate.)

If one has a whole *vector* of observations then the Cramer-Rao inequality involves the *joint* density function:

$$(10.4.29) \quad \text{var}[\mathbf{t}] \geq \frac{1}{E[(\frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta))^2]} = \frac{-1}{E[\frac{\partial^2}{\partial \theta^2} \log f_{\mathbf{y}}(\mathbf{y}; \theta)]}.$$

This inequality also holds if  $\mathbf{y}$  is discrete and one uses its probability mass function instead of the density function. In small samples, this lower bound is not always attainable; in some cases there is no unbiased estimator with a variance as low as the Cramer Rao lower bound.



PROBLEM 174. 4 points Assume  $n$  independent observations of a variable  $\mathbf{y} \sim \mathbf{N}(\mu, \sigma^2)$  are available, where  $\sigma^2$  is known. Show that the sample mean  $\bar{y}$  attains the Cramer-Rao lower bound for  $\mu$ .

ANSWER. The density function of each  $y_i$  is

$$(10.4.30) \quad f_{y_i}(y) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$$

therefore the log likelihood function of the whole vector is

$$(10.4.31) \quad \ell(\mathbf{y}; \mu) = \sum_{i=1}^n \log f_{y_i}(y_i) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2$$

$$(10.4.32) \quad \frac{\partial}{\partial \mu} \ell(\mathbf{y}; \mu) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu)$$

In order to apply (10.4.29) you can either square this and take the expected value

$$(10.4.33) \quad \mathbb{E}\left[\left(\frac{\partial}{\partial \mu} \ell(\mathbf{y}; \mu)\right)^2\right] = \frac{1}{\sigma^4} \sum \mathbb{E}[(y_i - \mu)^2] = n/\sigma^2$$

alternatively one may take one more derivative from (10.4.32) to get

$$(10.4.34) \quad \frac{\partial^2}{\partial \mu^2} \ell(\mathbf{y}; \mu) = -\frac{n}{\sigma^2}$$

This is constant, therefore equal to its expected value. Therefore the Cramer-Rao Lower Bound says that  $\text{var}[\bar{y}] \geq \sigma^2/n$ . This holds with equality.  $\square$

PROBLEM 175. Assume  $y_i \sim \text{NID}(0, \sigma^2)$  (i.e., normally independently distributed) with unknown  $\sigma^2$ . The obvious estimate of  $\sigma^2$  is  $s^2 = \frac{1}{n} \sum y_i^2$ .

• a. 2 points Show that  $s^2$  is an unbiased estimator of  $\sigma^2$ , is distributed  $\sim \frac{\sigma^2}{n} \chi_n^2$ , and has variance  $2\sigma^4/n$ . You are allowed to use the fact that a  $\chi_n^2$  has variance  $2n$ , which is equation (4.9.5).

ANSWER.

$$(10.4.35) \quad \mathbb{E}[y_i^2] = \text{var}[y_i] + (\mathbb{E}[y_i])^2 = \sigma^2 + 0 = \sigma^2$$

$$(10.4.36) \quad z_i = \frac{y_i}{\sigma} \sim \text{NID}(0, 1)$$

$$(10.4.37) \quad y_i = \sigma z_i$$

$$(10.4.38) \quad y_i^2 = \sigma^2 z_i^2$$

$$(10.4.39) \quad \sum_{i=1}^n y_i^2 = \sigma^2 \sum_{i=1}^n z_i^2 \sim \sigma^2 \chi_n^2$$

$$(10.4.40) \quad \frac{1}{n} \sum_{i=1}^n y_i^2 = \frac{\sigma^2}{n} \sum_{i=1}^n z_i^2 \sim \frac{\sigma^2}{n} \chi_n^2$$

$$(10.4.41) \quad \text{var}\left[\frac{1}{n} \sum_{i=1}^n y_i^2\right] = \frac{\sigma^4}{n^2} \text{var}[\chi_n^2] = \frac{\sigma^4}{n^2} 2n = \frac{2\sigma^4}{n}$$

• b. 4 points Show that this variance is at the same time the Cramer Rao lower bound.

ANSWER.

$$(10.4.42) \quad \ell(y, \sigma^2) = \log f_y(y; \sigma^2) = -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{y^2}{2\sigma^2}$$

$$(10.4.43) \quad \frac{\partial \log f_y}{\partial \sigma^2}(y; \sigma^2) = -\frac{1}{2\sigma^2} + \frac{y^2}{2\sigma^4} = \frac{y^2 - \sigma^2}{2\sigma^4}$$

Since  $\frac{y^2 - \sigma^2}{2\sigma^4}$  has zero mean, it follows

$$(10.4.44) \quad \mathbb{E}\left[\left(\frac{\partial \log f_y}{\partial \sigma^2}(y; \sigma^2)\right)^2\right] = \frac{\text{var}[y^2]}{4\sigma^8} = \frac{1}{2\sigma^4}.$$

Alternatively, one can differentiate one more time:

$$(10.4.45) \quad \frac{\partial^2 \log f_y}{(\partial \sigma^2)^2}(y; \sigma^2) = -\frac{y^2}{\sigma^6} + \frac{1}{2\sigma^4}$$

$$(10.4.46) \quad \mathbb{E}\left[\frac{\partial^2 \log f_y}{(\partial \sigma^2)^2}(y; \sigma^2)\right] = -\frac{\sigma^2}{\sigma^6} + \frac{1}{2\sigma^4} = \frac{1}{2\sigma^4}$$

$$(10.4.47)$$

This makes the Cramer Rao lower bound  $2\sigma^4/n$ .

PROBLEM 176. 4 points Assume  $x_1, \dots, x_n$  is a random sample of independent observations of a Poisson distribution with parameter  $\lambda$ , i.e., each of the  $x_i$  has probability mass function

$$(10.4.48) \quad p_{x_i}(x) = \Pr[x_i = x] = \frac{\lambda^x}{x!} e^{-\lambda} \quad x = 0, 1, 2, \dots$$

A Poisson variable with parameter  $\lambda$  has expected value  $\lambda$  and variance  $\lambda$ . (You are not required to prove this here.) Is there an unbiased estimator of  $\lambda$  with lower variance than the sample mean  $\bar{x}$ ?

Here is a formulation of the Cramer Rao Inequality for probability mass functions, as you need it for Question 176. Assume  $\mathbf{y}_1, \dots, \mathbf{y}_n$  are  $n$  independent observations of a random variable  $\mathbf{y}$  whose probability mass function depends on the unknown parameter  $\theta$  and satisfies certain regularity conditions. Write the univariate probability mass function of each of the  $\mathbf{y}_i$  as  $p_{\mathbf{y}}(y; \theta)$  and let  $\mathbf{t}$  be any unbiased estimator of  $\theta$ . Then

$$(10.4.49) \quad \text{var}[\mathbf{t}] \geq \frac{1}{n \mathbb{E}\left[\left(\frac{\partial}{\partial \theta} \ln p_{\mathbf{y}}(\mathbf{y}; \theta)\right)^2\right]} = \frac{-1}{n \mathbb{E}\left[\frac{\partial^2}{\partial \theta^2} \ln p_{\mathbf{y}}(\mathbf{y}; \theta)\right]}.$$

ANSWER. The Cramer Rao lower bound says no.

$$(10.4.50) \quad \log p_x(x; \lambda) = x \log \lambda - \log x! - \lambda$$

$$(10.4.51) \quad \frac{\partial \log p_x}{\partial \lambda}(x; \lambda) = \frac{x}{\lambda} - 1 = \frac{x - \lambda}{\lambda}$$

$$(10.4.52) \quad \mathbb{E}\left[\left(\frac{\partial \log p_x}{\partial \lambda}(x; \lambda)\right)^2\right] = \mathbb{E}\left[\frac{(x - \lambda)^2}{\lambda^2}\right] = \frac{\text{var}[x]}{\lambda^2} = \frac{1}{\lambda}.$$

Or alternatively, after (10.4.51) do

$$(10.4.53) \quad \frac{\partial^2 \log p_x}{\partial \lambda^2}(x; \lambda) = -\frac{x}{\lambda^2}$$

$$(10.4.54) \quad -\mathbb{E}\left[\frac{\partial^2 \log p_x}{\partial \lambda^2}(x; \lambda)\right] = \frac{\mathbb{E}[x]}{\lambda^2} = \frac{1}{\lambda}.$$

Therefore the Cramer Rao lower bound is  $\frac{\lambda}{n}$ , which is the variance of the sample mean.  $\square$

If the density function depends on more than one unknown parameter, i.e., if it has the form  $f_{\mathbf{y}}(y; \theta_1, \dots, \theta_k)$ , the Cramer Rao Inequality involves the following steps: (1) define  $\ell(\mathbf{y}; \theta_1, \dots, \theta_k) = \log f_{\mathbf{y}}(\mathbf{y}; \theta_1, \dots, \theta_k)$ , (2) form the following matrix which is called the *information matrix*:

$$(10.4.55) \quad \mathbf{I} = \begin{bmatrix} -n \mathbb{E}\left[\frac{\partial^2 \ell}{\partial \theta_1^2}\right] & \cdots & -n \mathbb{E}\left[\frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_k}\right] \\ \vdots & \ddots & \vdots \\ -n \mathbb{E}\left[\frac{\partial^2 \ell}{\partial \theta_k \partial \theta_1}\right] & \cdots & -n \mathbb{E}\left[\frac{\partial^2 \ell}{\partial \theta_k^2}\right] \end{bmatrix} = \begin{bmatrix} n \mathbb{E}\left[\left(\frac{\partial \ell}{\partial \theta_1}\right)^2\right] & \cdots & n \mathbb{E}\left[\frac{\partial \ell}{\partial \theta_1} \frac{\partial \ell}{\partial \theta_k}\right] \\ \vdots & \ddots & \vdots \\ n \mathbb{E}\left[\frac{\partial \ell}{\partial \theta_k} \frac{\partial \ell}{\partial \theta_1}\right] & \cdots & n \mathbb{E}\left[\left(\frac{\partial \ell}{\partial \theta_k}\right)^2\right] \end{bmatrix},$$

and (3) form the matrix inverse  $\mathbf{I}^{-1}$ . If the vector random variable  $\mathbf{t} = \begin{bmatrix} t_1 \\ \vdots \\ t_n \end{bmatrix}$

is an unbiased estimator of the parameter vector  $\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$ , then the inverse

the information matrix  $\mathbf{I}^{-1}$  is a lower bound for the covariance matrix  $\mathcal{V}[\mathbf{t}]$  in the following sense: the difference matrix  $\mathcal{V}[\mathbf{t}] - \mathbf{I}^{-1}$  is always nonnegative definite.

From this follows in particular: if  $i^{ii}$  is the  $i$ th diagonal element of  $\mathbf{I}^{-1}$ , then  $\text{var}[t_i] \geq i^{ii}$ .

### 10.5. Best Linear Unbiased Without Distribution Assumptions

If the  $x_i$  are Normal with unknown expected value and variance, their sample mean has lowest MSE among all unbiased estimators of  $\mu$ . If one does not assume Normality, then the sample mean has lowest MSE in the class of all *linear* unbiased estimators of  $\mu$ . This is true not only for the sample mean but also for all least squares estimates. This result needs remarkably weak assumptions: nothing is assumed about the distribution of the  $x_i$  other than the existence of mean and variance. Problem 177 shows that in some situations one can even dispense with the independence of the observations.

PROBLEM 177. 5 points [Lar82, example 5.4.1 on p 266] Let  $\mathbf{y}_1$  and  $\mathbf{y}_2$  be two random variables with same mean  $\mu$  and variance  $\sigma^2$ , but we do not assume that they are uncorrelated; their correlation coefficient is  $\rho$ , which can take any value  $|\rho| \leq 1$ . Show that  $\bar{\mathbf{y}} = (\mathbf{y}_1 + \mathbf{y}_2)/2$  has lowest mean squared error among all linear unbiased estimators of  $\mu$ , and compute its MSE. (An estimator  $\tilde{\mu}$  of  $\mu$  is linear iff it can be written in the form  $\tilde{\mu} = \alpha_1 \mathbf{y}_1 + \alpha_2 \mathbf{y}_2$  with some constant numbers  $\alpha_1$  and  $\alpha_2$ .)

ANSWER.

$$(10.5.1) \quad \bar{\mathbf{y}} = \alpha_1 \mathbf{y}_1 + \alpha_2 \mathbf{y}_2$$

$$(10.5.2) \quad \text{var} \bar{\mathbf{y}} = \alpha_1^2 \text{var}[\mathbf{y}_1] + \alpha_2^2 \text{var}[\mathbf{y}_2] + 2\alpha_1 \alpha_2 \text{cov}[\mathbf{y}_1, \mathbf{y}_2]$$

$$(10.5.3) \quad = \sigma^2(\alpha_1^2 + \alpha_2^2 + 2\alpha_1 \alpha_2 \rho).$$

Here we used (6.1.14). Unbiasedness means  $\alpha_2 = 1 - \alpha_1$ , therefore we call  $\alpha_1 = \alpha$  and  $\alpha_2 = 1 - \alpha$ .

$$(10.5.4) \quad \text{var}[\bar{\mathbf{y}}]/\sigma^2 = \alpha^2 + (1 - \alpha)^2 + 2\alpha(1 - \alpha)\rho$$

Now sort by the powers of  $\alpha$ :

$$(10.5.5) \quad = 2\alpha^2(1 - \rho) - 2\alpha(1 - \rho) + 1$$

$$(10.5.6) \quad = 2(\alpha^2 - \alpha)(1 - \rho) + 1.$$

This takes its minimum value where the derivative  $\frac{\partial}{\partial \alpha}(\alpha^2 - \alpha) = 2\alpha - 1 = 0$ . For the MSE plug  $\alpha_1 = \alpha_2 - 1/2$  into (10.5.3) to get  $\frac{\sigma^2}{2}(1 + \rho)$ .  $\square$

**PROBLEM 178.** *You have two unbiased measurements with errors of the same quantity  $\mu$  (which may or may not be random). The first measurement  $\mathbf{y}_1$  has mean squared error  $E[(\mathbf{y}_1 - \mu)^2] = \sigma^2$ , the other measurement  $\mathbf{y}_2$  has  $E[(\mathbf{y}_1 - \mu)^2] = \tau^2$ . The measurement errors  $\mathbf{y}_1 - \mu$  and  $\mathbf{y}_2 - \mu$  have zero expected values (i.e., the measurements are unbiased) and are independent of each other.*

• a. 2 points Show that the linear unbiased estimators of  $\mu$  based on these two measurements are simply the weighted averages of these measurements, i.e., they can be written in the form  $\tilde{\mu} = \alpha\mathbf{y}_1 + (1 - \alpha)\mathbf{y}_2$ , and that the MSE of such an estimator is  $\alpha^2\sigma^2 + (1 - \alpha)^2\tau^2$ . Note: we are using the word “estimator” here even if  $\mu$  is random. An estimator or predictor  $\tilde{\mu}$  is unbiased if  $E[\tilde{\mu} - \mu] = 0$ . Since we allow  $\mu$  to be random, the proof in the class notes has to be modified.

**ANSWER.** The estimator  $\tilde{\mu}$  is linear (more precisely: affine) if it can be written in the form

$$(10.5.7) \quad \tilde{\mu} = \alpha_1\mathbf{y}_1 + \alpha_2\mathbf{y}_2 + \gamma$$

The measurements themselves are unbiased, i.e.,  $E[\mathbf{y}_i - \mu] = 0$ , therefore

$$(10.5.8) \quad E[\tilde{\mu} - \mu] = (\alpha_1 + \alpha_2 - 1)E[\mu] + \gamma = 0$$

for all possible values of  $E[\mu]$ ; therefore  $\gamma = 0$  and  $\alpha_2 = 1 - \alpha_1$ . To simplify notation, we will call from now on  $\alpha_1 = \alpha$ ,  $\alpha_2 = 1 - \alpha$ . Due to unbiasedness, the MSE is the variance of the estimation error

$$(10.5.9) \quad \text{var}[\tilde{\mu} - \mu] = \alpha^2\sigma^2 + (1 - \alpha)^2\tau^2$$

$\square$

• b. 4 points Define  $\omega^2$  by

$$(10.5.10) \quad \frac{1}{\omega^2} = \frac{1}{\sigma^2} + \frac{1}{\tau^2} \quad \text{which can be solved to give} \quad \omega^2 = \frac{\sigma^2\tau^2}{\sigma^2 + \tau^2}.$$

Show that the Best (i.e., minimum MSE) linear unbiased estimator (BLUE) of  $\mu$  based on these two measurements is

$$(10.5.11) \quad \hat{\mathbf{y}} = \frac{\omega^2}{\sigma^2}\mathbf{y}_1 + \frac{\omega^2}{\tau^2}\mathbf{y}_2$$

i.e., it is the weighted average of  $\mathbf{y}_1$  and  $\mathbf{y}_2$  where the weights are proportional to the inverses of the variances.

**ANSWER.** The variance (10.5.9) takes its minimum value where its derivative with respect to  $\alpha$  is zero, i.e., where

$$(10.5.12) \quad \frac{\partial}{\partial \alpha}(\alpha^2\sigma^2 + (1 - \alpha)^2\tau^2) = 2\alpha\sigma^2 - 2(1 - \alpha)\tau^2 = 0$$

$$(10.5.13) \quad \alpha\sigma^2 = \tau^2 - \alpha\tau^2$$

$$(10.5.14) \quad \alpha = \frac{\tau^2}{\sigma^2 + \tau^2}$$

In terms of  $\omega$  one can write

$$(10.5.15) \quad \alpha = \frac{\tau^2}{\sigma^2 + \tau^2} = \frac{\omega^2}{\sigma^2} \quad \text{and} \quad 1 - \alpha = \frac{\sigma^2}{\sigma^2 + \tau^2} = \frac{\omega^2}{\tau^2}.$$

• c. 2 points Show: the MSE of the BLUE  $\omega^2$  satisfies the following equation

$$(10.5.16) \quad \frac{1}{\omega^2} = \frac{1}{\sigma^2} + \frac{1}{\tau^2}$$

**ANSWER.** We already have introduced the notation  $\omega^2$  for the quantity defined by (10.5.10); therefore all we have to show is that the MSE or, equivalently, the variance of the estimation error is equal to this  $\omega^2$ :

$$(10.5.17) \quad \text{var}[\tilde{\mu} - \mu] = \left(\frac{\omega^2}{\sigma^2}\right)^2\sigma^2 + \left(\frac{\omega^2}{\tau^2}\right)^2\tau^2 = \omega^4\left(\frac{1}{\sigma^2} + \frac{1}{\tau^2}\right) = \omega^4\frac{1}{\omega^2} = \omega^2$$

Examples of other classes of estimators for which a best estimator exists are one requires the estimator to be translation invariant, then the least squares estimators are best in the class of all translation invariant estimators. But there is no best linear estimator in the linear model. (Theil)

## 10.6. Maximum Likelihood Estimation

This is an excellent and very widely applicable estimation principle. Its major drawback is its computational complexity, but with modern computing power it becomes more and more manageable. Another drawback is that it requires a full specification of the distribution.

**PROBLEM 179.** 2 points What are the two greatest disadvantages of Maximum Likelihood Estimation?

**ANSWER.** Its high information requirements (the functional form of the density function must be known), and computational complexity.

In our discussion of entropy in Section 3.11 we derived an extremal value property which distinguishes the actual density function  $f_{\mathbf{y}}(y)$  of a given random variable from all other possible density functions of  $\mathbf{y}$ , i.e., from all other functions  $g \geq 0$  with  $\int_{-\infty}^{+\infty} g(y) dy = 1$ . The true density function of  $\mathbf{y}$  is the one which maximizes

$E[\log g(\mathbf{y})]$ . We showed that this principle can be used to design a payoff scheme by which it is in the best interest of a forecaster to tell the truth. Now we will see that this principle can also be used to design a good estimator. Say you have  $n$  independent observations of  $\mathbf{y}$ . You know the density of  $\mathbf{y}$  belongs to a given family  $\mathcal{F}$  of density functions, but you don't know which member of  $\mathcal{F}$  it is. Then form the arithmetic mean of  $\log f(y_i)$  for all  $f \in \mathcal{F}$ . It converges towards  $E[\log f(\mathbf{y})]$ . For the true density function, this expected value is higher than for all the other density functions. If one does not know which the true density function is, then it is a good strategy to select that density function  $f$  for which the sample mean of the  $\log f(y_i)$  is largest. This is the maximum likelihood estimator.

Let us interject here a short note about the definitional difference between density function and likelihood function. If we know  $\mu = \mu_0$ , we can write down the density function as

$$(10.6.1) \quad f_{\mathbf{y}}(y; \mu_0) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu_0)^2}{2}}$$

It is a function of  $y$ , the possible values assumed by  $\mathbf{y}$ , and the letter  $\mu_0$  symbolizes a constant, the true parameter value. The same function considered as a function of the *variable*  $\mu$ , representing all possible values assumable by the true mean, with  $y$  being *fixed* at the actually observed value, becomes the likelihood function.

In the same way one can also turn probability mass functions  $p_{\mathbf{x}}(x)$  into likelihood functions.

Now let us compute some examples of the MLE. You make  $n$  independent observations  $\mathbf{y}_1, \dots, \mathbf{y}_n$  from a  $N(\mu, \sigma^2)$  distribution. Write the likelihood function as

$$(10.6.2) \quad L(\mu, \sigma^2; \mathbf{y}_1, \dots, \mathbf{y}_n) = \prod_{i=1}^n f_{\mathbf{y}}(\mathbf{y}_i) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n e^{-\frac{1}{2\sigma^2} \sum (\mathbf{y}_i - \mu)^2}.$$

Its logarithm is more convenient to maximize:

$$(10.6.3) \quad \ell = \ln L(\mu, \sigma^2; \mathbf{y}_1, \dots, \mathbf{y}_n) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum (\mathbf{y}_i - \mu)^2.$$

To compute the maximum we need the partial derivatives:

$$(10.6.4) \quad \frac{\partial \ell}{\partial \mu} = \frac{1}{\sigma^2} \sum (\mathbf{y}_i - \mu)$$

$$(10.6.5) \quad \frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (\mathbf{y}_i - \mu)^2.$$

The maximum likelihood estimators are those values  $\hat{\mu}$  and  $\hat{\sigma}^2$  which set these two partials zero. I.e., at the same time at which we set the partials zero we must put the hats on  $\mu$  and  $\sigma^2$ . As long as  $\hat{\sigma}^2 \neq 0$  (which is the case with probability one), the first equation determines  $\hat{\mu}$ :  $\sum \mathbf{y}_i - n\hat{\mu} = 0$ , i.e.,  $\hat{\mu} = \frac{1}{n} \sum \mathbf{y}_i = \bar{\mathbf{y}}$ . (This would

be the MLE of  $\mu$  even if  $\sigma^2$  were known). Now plug this  $\hat{\mu}$  into the second equation to get  $\frac{n}{2} = \frac{1}{2\hat{\sigma}^2} \sum (\mathbf{y}_i - \bar{\mathbf{y}})^2$ , or  $\hat{\sigma}^2 = \frac{1}{n} \sum (\mathbf{y}_i - \bar{\mathbf{y}})^2$ .

Here is another example:  $\mathbf{t}_1, \dots, \mathbf{t}_n$  are independent and follow an exponential distribution, i.e.,

$$(10.6.6) \quad f_{\mathbf{t}}(t; \lambda) = \lambda e^{-\lambda t} \quad (t > 0)$$

$$(10.6.7) \quad L(\mathbf{t}_1, \dots, \mathbf{t}_n; \lambda) = \lambda^n e^{-\lambda(\mathbf{t}_1 + \dots + \mathbf{t}_n)}$$

$$(10.6.8) \quad \ell(\mathbf{t}_1, \dots, \mathbf{t}_n; \lambda) = n \ln \lambda - \lambda(\mathbf{t}_1 + \dots + \mathbf{t}_n)$$

$$(10.6.9) \quad \frac{\partial \ell}{\partial \lambda} = \frac{n}{\lambda} - (\mathbf{t}_1 + \dots + \mathbf{t}_n).$$

Set this zero, and write  $\hat{\lambda}$  instead of  $\lambda$  to get  $\hat{\lambda} = \frac{n}{\mathbf{t}_1 + \dots + \mathbf{t}_n} = 1/\bar{\mathbf{t}}$ .

Usually the MLE is asymptotically unbiased and asymptotically normal. Therefore it is important to have an estimate of its asymptotic variance. Here we can use the fact that asymptotically the Cramer Rao Lower Bound is not merely a lower bound for this variance but is equal to its variance. (From this follows that the maximum likelihood estimator is asymptotically efficient.) The Cramer Rao lower bound itself depends on unknown parameters. In order to get a consistent estimate of the Cramer Rao lower bound, do the following: (1) Replace the unknown parameters in the second derivative of the log likelihood function by their maximum likelihood estimates. (2) Instead of taking expected values over the observed values  $\mathbf{x}_i$  you may simply insert the sample values of the  $\mathbf{x}_i$  into these maximum likelihood estimates and (3) then invert this estimate of the information matrix.

MLE obeys an important functional invariance principle: if  $\hat{\theta}$  is the MLE of  $\theta$  then  $g(\hat{\theta})$  is the MLE of  $g(\theta)$ . E.g.,  $\mu = \frac{1}{\lambda}$  is the expected value of the exponential variable, and its MLE is  $\bar{\mathbf{x}}$ .

**PROBLEM 180.**  $\mathbf{x}_1, \dots, \mathbf{x}_m$  is a sample from a  $N(\mu_{\mathbf{x}}, \sigma^2)$ , and  $\mathbf{y}_1, \dots, \mathbf{y}_n$  from a  $N(\mu_{\mathbf{y}}, \sigma^2)$  with different mean but same  $\sigma^2$ . All observations are independent of each other.

• a. 2 points Show that the MLE of  $\mu_{\mathbf{x}}$ , based on the combined sample, is  $\bar{\mathbf{x}}$ . (By symmetry it follows that the MLE of  $\mu_{\mathbf{y}}$  is  $\bar{\mathbf{y}}$ .)

ANSWER.

$$(10.6.10) \quad \ell(\mu_{\mathbf{x}}, \mu_{\mathbf{y}}, \sigma^2) = -\frac{m}{2} \ln 2\pi - \frac{m}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^m (\mathbf{x}_i - \mu_{\mathbf{x}})^2$$

$$- \frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^n (\mathbf{y}_j - \mu_{\mathbf{y}})^2$$

$$(10.6.11) \quad \frac{\partial \ell}{\partial \mu_{\mathbf{x}}} = -\frac{1}{2\sigma^2} \sum -2(\mathbf{x}_i - \mu_{\mathbf{x}}) = 0 \quad \text{for } \mu_{\mathbf{x}} = \bar{\mathbf{x}}$$

- b. 2 points Derive the MLE of  $\sigma^2$ , based on the combined samples.

ANSWER.

$$(10.6.12) \quad \frac{\partial \ell}{\partial \sigma^2} = -\frac{m+n}{2\sigma^2} + \frac{1}{2\sigma^4} \left( \sum_{i=1}^m (x_i - \mu_x)^2 + \sum_{j=1}^n (y_j - \mu_y)^2 \right)$$

$$(10.6.13) \quad \hat{\sigma}^2 = \frac{1}{m+n} \left( \sum_{i=1}^m (x_i - \bar{x})^2 + \sum_{j=1}^n (y_j - \bar{y})^2 \right).$$

□

### 10.7. Method of Moments Estimators

Method of moments estimators use the sample moments as estimates of the population moments. I.e., the estimate of  $\mu$  is  $\bar{x}$ , the estimate of the variance  $\sigma^2$  is  $\frac{1}{n} \sum (x_i - \bar{x})^2$ , etc. If the parameters are a given function of the population moments, use the same function of the sample moments (using the lowest moments which do the job).

The advantage of method of moments estimators is their computational simplicity. Many of the estimators discussed above are method of moments estimators. However if the moments do not exist, then method of moments estimators are inconsistent, and in general method of moments estimators are not as good as maximum likelihood estimators.

### 10.8. M-Estimators

The class of  $M$ -estimators maximizes something other than a likelihood function: it includes nonlinear least squares, generalized method of moments, minimum distance and minimum chi-squared estimators. The purpose is to get a “robust” estimator which is good for a wide variety of likelihood functions. Many of these are asymptotically efficient; but their small-sample properties may vary greatly.

### 10.9. Sufficient Statistics and Estimation

*Weak Sufficiency Principle:* If  $\mathbf{x}$  has a p.d.f.  $f_{\mathbf{x}}(\mathbf{x}; \theta)$  and if a sufficient statistic  $\mathbf{s}(\mathbf{x})$  exists for  $\theta$ , then identical conclusions should be drawn from data  $\mathbf{x}_1$  and  $\mathbf{x}_2$  which have same value  $\mathbf{s}(\mathbf{x}_1) = \mathbf{s}(\mathbf{x}_2)$ .

Why? Sufficiency means: after knowing  $\mathbf{s}(\mathbf{x})$ , the rest of the data  $\mathbf{x}$  can be regarded generated by a random mechanism not dependent on  $\theta$ , and are therefore uninformative about  $\theta$ .

This principle can be used to improve on given estimators. Without proof we will state here

*Rao Blackwell Theorem:* Let  $\mathbf{t}(\mathbf{x})$  be an estimator of  $\theta$  and  $\mathbf{s}(\mathbf{x})$  a sufficient statistic for  $\theta$ . Then one can get an estimator  $\mathbf{t}^*(\mathbf{x})$  of  $\theta$  which has no worse MSE than  $\mathbf{t}(\mathbf{x})$  by taking expectations conditionally on the sufficient statistic, i.e.  $\mathbf{t}^*(\mathbf{x}) = E[\mathbf{t}(\mathbf{x})|\mathbf{s}(\mathbf{x})]$ .

To recapitulate:  $\mathbf{t}^*(\mathbf{x})$  is obtained by the following two steps: (1) Compute the conditional expectation  $\mathbf{t}^{**}(\mathbf{s}) = E[\mathbf{t}(\mathbf{x})|\mathbf{s}(\mathbf{x}) = \mathbf{s}]$ , and (2) plug  $\mathbf{s}(\mathbf{x})$  into  $\mathbf{t}^{**}$ , i.e.  $\mathbf{t}^*(\mathbf{x}) = \mathbf{t}^{**}(\mathbf{s}(\mathbf{x}))$ .

A statistic  $\mathbf{s}$  is said to be *complete*, if the only real-valued function  $g$  defined on the range of  $\mathbf{s}$ , which satisfies  $E[g(\mathbf{s})] = 0$  whatever the value of  $\theta$ , is the function which is identically zero. If a statistic  $\mathbf{s}$  is complete and sufficient, then every function  $g(\mathbf{s})$  is the minimum MSE unbiased estimator of its expected value  $E[g(\mathbf{s})]$ .

If a complete and sufficient statistic exists, this gives a systematic approach to minimum MSE unbiased estimators (*Lehmann-Scheffé Theorem*): if  $\mathbf{t}$  is an unbiased estimator of  $\theta$  and  $\mathbf{s}$  is *complete* and sufficient, then  $\mathbf{t}^*(\mathbf{x}) = E[\mathbf{t}(\mathbf{x})|\mathbf{s}(\mathbf{x})]$  has low MSE in the class of all unbiased estimators of  $\theta$ . Problem 181 steps you through the proof.

PROBLEM 181. [BD77, Problem 4.2.6 on p. 144] *If a statistic  $\mathbf{s}$  is complete and sufficient, then every function  $g(\mathbf{s})$  is the minimum MSE unbiased estimator of  $E[g(\mathbf{s})]$  (Lehmann-Scheffé theorem). This gives a systematic approach to find minimum MSE unbiased estimators. Here are the definitions:  $\mathbf{s}$  is sufficient for  $\theta$  if for any event  $E$  and any value  $\mathbf{s}$ , the conditional probability  $\Pr[E|\mathbf{s} \leq \mathbf{s}]$  does not involve  $\theta$ .  $\mathbf{s}$  is complete for  $\theta$  if the only function  $g(\mathbf{s})$  of  $\mathbf{s}$ , which has zero expected value whatever the value of  $\theta$ , is the function which is identically zero, i.e.,  $g(\mathbf{s}) = 0$  for all  $\mathbf{s}$ .*

- a. 3 points Given an unknown parameter  $\theta$ , and a complete sufficient statistic  $\mathbf{s}$ , how can one find that function of  $\mathbf{s}$  whose expected value is  $\theta$ ? There is an easy trick: start with any statistic  $\mathbf{p}$  with  $E[\mathbf{p}] = \theta$ , and use the conditional expectation  $E[\mathbf{p}|\mathbf{s}]$ . Argue why this conditional expectation does not depend on the unknown parameter  $\theta$ , is an unbiased estimator of  $\theta$ , and why this leads to the same estimator regardless which  $\mathbf{p}$  one starts with.

ANSWER. You need sufficiency for the first part of the problem, the law of iterated expectations for the second, and completeness for the third.

Set  $E = \{\mathbf{p} \leq \mathbf{p}\}$  in the definition of sufficiency given at the beginning of the Problem to see that the cdf of  $\mathbf{p}$  conditionally on  $\mathbf{s}$  being in any interval does not involve  $\theta$ , therefore also  $E[\mathbf{p}|\mathbf{s}]$  does not involve  $\theta$ .

Unbiasedness follows from the theorem of iterated expectations  $E[E[\mathbf{p}|\mathbf{s}]] = E[\mathbf{p}] = \theta$ .

The independence on the choice of  $\mathbf{p}$  can be shown as follows: Since the conditional expectation conditionally on  $\mathbf{s}$  is a function of  $\mathbf{s}$ , we can use the notation  $E[\mathbf{p}|\mathbf{s}] = g_1(\mathbf{s})$  and  $E[\mathbf{q}|\mathbf{s}] = g_2(\mathbf{s})$ . From  $E[\mathbf{p}] = E[\mathbf{q}]$  follows by the law of iterated expectations  $E[g_1(\mathbf{s}) - g_2(\mathbf{s})] = 0$ , therefore by completeness  $g_1(\mathbf{s}) - g_2(\mathbf{s}) \equiv 0$ .

• b. 2 points Assume  $\mathbf{y}_i \sim \text{NID}(\mu, 1)$  ( $i = 1, \dots, n$ ), i.e., they are independent and normally distributed with mean  $\mu$  and variance 1. Without proof you are allowed to use the fact that in this case, the sample mean  $\bar{\mathbf{y}}$  is a complete sufficient statistic for  $\mu$ . What is the minimum MSE unbiased estimate of  $\mu$ , and what is that of  $\mu^2$ ?

ANSWER. We have to find functions of  $\bar{\mathbf{y}}$  with the desired parameters as expected values. Clearly,  $\bar{\mathbf{y}}$  is that of  $\mu$ , and  $\bar{\mathbf{y}}^2 - 1/n$  is that of  $\mu^2$ .  $\square$

• c. 1 point For a given  $j$ , let  $\pi$  be the probability that the  $j^{\text{th}}$  observation is nonnegative, i.e.,  $\pi = \Pr[\mathbf{y}_j \geq 0]$ . Show that  $\pi = \Phi(\mu)$  where  $\Phi$  is the cumulative distribution function of the standard normal. The purpose of the remainder of this Problem is to find a minimum MSE unbiased estimator of  $\pi$ .

ANSWER.

$$(10.9.1) \quad \pi = \Pr[\mathbf{y}_i \geq 0] = \Pr[\mathbf{y}_i - \mu \geq -\mu] = \Pr[\mathbf{y}_i - \mu \leq \mu] = \Phi(\mu)$$

because  $\mathbf{y}_i - \mu \sim \text{N}(0, 1)$ . We needed symmetry of the distribution to flip the sign.  $\square$

• d. 1 point As a first step we have to find an unbiased estimator of  $\pi$ . It does not have to be a good one, any unbiased estimator will do. And such an estimator is indeed implicit in the definition of  $\pi$ . Let  $q$  be the “indicator function” for nonnegative values, satisfying  $q(y) = 1$  if  $y \geq 0$  and 0 otherwise. We will be working with the random variable which one obtains by inserting the  $j^{\text{th}}$  observation  $\mathbf{y}_j$  into  $q$ , i.e., with  $\mathbf{q} = q(\mathbf{y}_j)$ . Show that  $\mathbf{q}$  is an unbiased estimator of  $\pi$ .

ANSWER.  $q(\mathbf{y}_j)$  has a discrete distribution and  $\Pr[q(\mathbf{y}_j) = 1] = \Pr[\mathbf{y}_j \geq 0] = \pi$  by (10.9.1) and therefore  $\Pr[q(\mathbf{y}_j) = 0] = 1 - \pi$

The expected value is  $E[q(\mathbf{y}_j)] = (1 - \pi) \cdot 0 + \pi \cdot 1 = \pi$ .  $\square$

• e. 2 points Given  $\mathbf{q}$  we can apply the Lehmann-Scheffé theorem:  $E[q(\mathbf{y}_j)|\bar{\mathbf{y}}]$  is the best unbiased estimator of  $\pi$ . We will compute  $E[q(\mathbf{y}_j)|\bar{\mathbf{y}}]$  in four steps which build on each other. First step: since for every indicator function follows  $E[q(\mathbf{y}_j)|\bar{\mathbf{y}}] = \Pr[\mathbf{y}_j \geq 0|\bar{\mathbf{y}}]$ , we need for every given value  $\bar{\mathbf{y}}$ , the conditional distribution of  $\mathbf{y}_j$  conditionally on  $\bar{\mathbf{y}} = \bar{\mathbf{y}}$ . (Not just the conditional mean but the whole conditional distribution.) In order to construct this, we first have to specify exactly the joint distribution of  $\mathbf{y}_j$  and  $\bar{\mathbf{y}}$ :

ANSWER. They are jointly normal:

$$(10.9.2) \quad \begin{bmatrix} \mathbf{y}_j \\ \bar{\mathbf{y}} \end{bmatrix} \sim N \left( \begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} 1 & 1/n \\ 1/n & 1/n \end{bmatrix} \right)$$

$\square$

• f. 2 points Second step: From this joint distribution derive the conditional distribution of  $\mathbf{y}_j$  conditionally on  $\bar{\mathbf{y}} = \bar{\mathbf{y}}$ . (Not just the conditional mean but the whole conditional distribution.) For this you will need formula (7.3.18) and (7.3.20).

ANSWER. Here are these two formulas: if  $\mathbf{u}$  and  $\mathbf{v}$  are jointly normal, then the conditional distribution of  $\mathbf{v}$  conditionally on  $\mathbf{u} = u$  is Normal with mean

$$(10.9.3) \quad E[\mathbf{v}|\mathbf{u} = u] = E[\mathbf{v}] + \frac{\text{cov}[\mathbf{u}, \mathbf{v}]}{\text{var}[\mathbf{u}]}(u - E[\mathbf{u}])$$

and variance

$$(10.9.4) \quad \text{var}[\mathbf{v}|\mathbf{u} = u] = \text{var}[\mathbf{v}] - \frac{(\text{cov}[\mathbf{u}, \mathbf{v}])^2}{\text{var}[\mathbf{u}]}.$$

Plugging  $\mathbf{u} = \bar{\mathbf{y}}$  and  $\mathbf{v} = \mathbf{y}_j$  into (7.3.18) and (7.3.20) gives: the conditional distribution of  $\mathbf{y}_j$  conditionally on  $\bar{\mathbf{y}} = \bar{\mathbf{y}}$  has mean

$$(10.9.5) \quad E[\mathbf{y}_j|\bar{\mathbf{y}} = \bar{\mathbf{y}}] = E[\mathbf{y}_j] + \frac{\text{cov}[\bar{\mathbf{y}}, \mathbf{y}_j]}{\text{var}[\bar{\mathbf{y}}]}(\bar{\mathbf{y}} - E[\bar{\mathbf{y}}])$$

$$(10.9.6) \quad = \mu + \frac{1/n}{1/n}(\bar{\mathbf{y}} - \mu) = \bar{\mathbf{y}}$$

and variance

$$(10.9.7) \quad \text{var}[\mathbf{y}_j|\bar{\mathbf{y}} = \bar{\mathbf{y}}] = \text{var}[\mathbf{y}_j] - \frac{(\text{cov}[\bar{\mathbf{y}}, \mathbf{y}_j])^2}{\text{var}[\bar{\mathbf{y}}]}$$

$$(10.9.8) \quad = 1 - \frac{(1/n)^2}{1/n} = 1 - \frac{1}{n}.$$

Therefore the conditional distribution of  $\mathbf{y}_j$  conditional on  $\bar{\mathbf{y}}$  is  $\text{N}(\bar{\mathbf{y}}, (n-1)/n)$ . How can this be motivated? if we know the actual arithmetic mean of the variables, then our best estimator that each variable is equal to this arithmetic mean. And this additional knowledge cuts down the variance by  $1/n$ .

• g. 2 points The variance decomposition (6.6.6) gives a decomposition of  $\text{var}[\mathbf{y}_j]$  give it here:

ANSWER.

$$(10.9.9) \quad \text{var}[\mathbf{y}_j] = \text{var}[E[\mathbf{y}_j|\bar{\mathbf{y}}]] + E[\text{var}[\mathbf{y}_j|\bar{\mathbf{y}}]]$$

$$(10.9.10) \quad = \text{var}[\bar{\mathbf{y}}] + E\left[\frac{n-1}{n}\right] = \frac{1}{n} + \frac{n-1}{n}$$

• h. Compare the conditional with the unconditional distribution.

ANSWER. Conditional distribution does not depend on unknown parameters, and it has same variance!

• i. 2 points Third step: Compute the probability, conditionally on  $\bar{\mathbf{y}} = \bar{\mathbf{y}}$ , that  $\mathbf{y}_j \geq 0$ .

ANSWER. If  $\mathbf{x} \sim \text{N}(\bar{\mathbf{y}}, (n-1)/n)$  (I call it  $\mathbf{x}$  here instead of  $\mathbf{y}_j$  since we use it not with its family unconditional distribution  $\text{N}(\mu, 1)$  but with a conditional distribution), then  $\Pr[\mathbf{x} \geq 0] = \Pr[\mathbf{x} - \bar{\mathbf{y}} \geq -\bar{\mathbf{y}}] = \Pr[\mathbf{x} - \bar{\mathbf{y}} \leq \bar{\mathbf{y}}] = \Pr[(\mathbf{x} - \bar{\mathbf{y}})\sqrt{n/(n-1)} \leq \bar{\mathbf{y}}\sqrt{n/(n-1)}] = \Phi(\bar{\mathbf{y}}\sqrt{n/(n-1)})$  because  $(\mathbf{x} - \bar{\mathbf{y}})\sqrt{n/(n-1)} \sim \text{N}(0, 1)$  conditionally on  $\bar{\mathbf{y}}$ . Again we needed symmetry of the distribution to flip the sign.

• j. 1 point Finally, put all the pieces together and write down  $E[q(\mathbf{y}_j)|\bar{\mathbf{y}}]$ , the conditional expectation of  $q(\mathbf{y}_j)$  conditionally on  $\bar{\mathbf{y}}$ , which by the Lehmann-Scheffé theorem is the minimum MSE unbiased estimator of  $\pi$ . The formula you should come up with is

$$(10.9.11) \quad \hat{\pi} = \Phi(\bar{\mathbf{y}}\sqrt{n/(n-1)}),$$

where  $\Phi$  is the standard normal cumulative distribution function.

ANSWER. The conditional expectation of  $q(\mathbf{y}_j)$  conditionally on  $\bar{\mathbf{y}} = \bar{y}$  is, by part d, simply the probability that  $\mathbf{y}_j \geq 0$  under this conditional distribution. In part i this was computed as  $\Phi(\bar{y}\sqrt{n/(n-1)})$ . Therefore all we have to do is replace  $\bar{y}$  by  $\bar{\mathbf{y}}$  to get the minimum MSE unbiased estimator of  $\pi$  as  $\Phi(\bar{\mathbf{y}}\sqrt{n/(n-1)})$ .  $\square$

*Remark: this particular example did not give any brand new estimators, but it can rather be considered a proof that certain obvious estimators are unbiased and efficient. But often this same procedure gives new estimators which one would not have been able to guess. Already when the variance is unknown, the above example becomes quite a bit more complicated, see [Rao73, p. 322, example 2]. When the variables have an exponential distribution then this example (probability of early failure) is discussed in [BD77, example 4.2.4 on pp. 124/5].*

### 10.10. The Likelihood Principle

Consider two experiments whose likelihood functions depend on the same parameter vector  $\theta$ . Suppose that for particular realizations of the data  $\mathbf{y}_1$  and  $\mathbf{y}_2$ , the respective likelihood functions are proportional to each other, i.e.,  $l_1(\theta; \mathbf{y}_1) = \alpha l_2(\theta; \mathbf{y}_2)$  where  $\alpha$  does not depend on  $\theta$  although it may depend on  $\mathbf{y}_1$  and  $\mathbf{y}_2$ . Then the *likelihood principle* states that identical conclusions should be drawn from these two experiments about  $\theta$ .

The likelihood principle is equivalent to the combination of two simpler principles: the weak sufficiency principle, and the following principle, which seems very plausible:

*Weak Conditionality Principle:* Given two possible experiments  $A$  and  $B$ . A mixed experiment is one in which one throws a coin and performs  $A$  if the coin shows head and  $B$  if it shows tails. The weak conditionality principle states: suppose it is known that the coin shows tails. Then the evidence of the mixed experiment is equivalent to the evidence gained had one not thrown the coin but performed  $B$  without the possible alternative of  $A$ . This principle says therefore that an experiment which one did not do but which one could have performed does not alter the information gained from the experiment actually performed.

As an application of the likelihood principle look at the following situation:

PROBLEM 182. 3 points You have a Bernoulli experiment with unknown parameter  $\theta$ ,  $0 \leq \theta \leq 1$ . Person  $A$  was originally planning to perform this experiment 12 times, which she does. She obtains 9 successes and 3 failures. Person  $B$  was originally planning to perform the experiment until he has reached 9 successes, and it took him 12 trials to do this. Should both experimenters draw identical conclusions from these two experiments or not?

ANSWER. The probability mass function in the first is by (3.7.1)  $\binom{12}{9}\theta^9(1-\theta)^3$ , and in second it is by (4.1.13)  $\binom{11}{8}\theta^9(1-\theta)^3$ . They are proportional, the stopping rule therefore does not matter!

### 10.11. Bayesian Inference

Real-life estimation usually implies the choice between competing estimation methods all of which have their advantages and disadvantages. Bayesian inference removes some of this arbitrariness.

Bayesians claim that “any inferential or decision process that does not follow from some likelihood function and some set of priors has objectively verifiable deficiencies” [Cor69, p. 617]. The “prior information” used by Bayesians is a formalization of the notion that the information about the parameter values never comes from the experiment alone. The Bayesian approach to estimation forces the researcher to combine his or her prior knowledge (and also the loss function for estimation errors) in a mathematical form, because in this way, unambiguous mathematical prescriptions can be derived as to how the information of an experiment should be evaluated.

To the objection that these are large information requirements which are often not satisfied, one might answer that it is less important whether these assumptions are actually the right ones. The formulation of prior density merely ensures that the researcher proceeds from a coherent set of beliefs.

The mathematics which the Bayesians do is based on a “final” instead of an “initial” criterion of precision. In other words, not an estimation procedure is evaluated which will be good in hypothetical repetitions of the experiment in the average, but one which is good for the given set of data and the given set of priors. Data which could have been observed but were not observed are not taken into consideration.

Both Bayesians and non-Bayesians define the probabilistic properties of an experiment by the density function (likelihood function) of the observations, which may depend on one or several unknown parameters. The non-Bayesian considers the parameters fixed but unknown, while the Bayesian considers the parameters random, i.e., he symbolizes his prior information about the parameters by a prior probability distribution.

An excellent example in which this prior probability distribution is discrete is given in [Ame94, pp. 168–172]. In the more usual case that the prior distribution

has a density function, a Bayesian is working with the joint density function of the parameter values and the data. Like all joint density function, it can be written as the product of a marginal and conditional density. The marginal density of the parameter value represents the beliefs the experimenter holds about the parameters before the experiment (prior density), and the likelihood function of the experiment is the conditional density of the data given the parameters. After the experiment has been conducted, the experimenter's belief about the parameter values is represented by their *conditional density given the data*, called the posterior density.

Let  $\mathbf{y}$  denote the observations,  $\boldsymbol{\theta}$  the unknown parameters, and  $f(\mathbf{y}, \boldsymbol{\theta})$  their joint density. Then

$$(10.11.1) \quad f(\mathbf{y}, \boldsymbol{\theta}) = f(\boldsymbol{\theta})f(\mathbf{y}|\boldsymbol{\theta})$$

$$(10.11.2) \quad = f(\mathbf{y})f(\boldsymbol{\theta}|\mathbf{y}).$$

Therefore

$$(10.11.3) \quad f(\boldsymbol{\theta}|\mathbf{y}) = \frac{f(\boldsymbol{\theta})f(\mathbf{y}|\boldsymbol{\theta})}{f(\mathbf{y})}.$$

In this formula, the value of  $f(\mathbf{y})$  is irrelevant. It only depends on  $\mathbf{y}$  but not on  $\boldsymbol{\theta}$ , but  $\mathbf{y}$  is fixed, i.e., it is a constant. If one knows the posterior density function of  $\boldsymbol{\theta}$  up to a constant, one knows it altogether, since the constant is determined by the requirement that the area under the density function is 1. Therefore (10.11.3) is usually written as ( $\propto$  means “proportional to”)

$$(10.11.4) \quad f(\boldsymbol{\theta}|\mathbf{y}) \propto f(\boldsymbol{\theta})f(\mathbf{y}|\boldsymbol{\theta});$$

here the lefthand side contains the posterior density function of the parameter, the righthand side the prior density function and the likelihood function representing the probability distribution of the experimental data.

The Bayesian procedure does not yield a point estimate or an interval estimate, but a whole probability distribution for the unknown parameters (which represents our *information about* these parameters) containing the “prior” information “updated” by the information yielded by the sample outcome.

Of course, such probability distributions can be summarized by various measures of location (mean, median), which can then be considered Bayesian point estimates. Such summary measures for a whole probability distribution are rather arbitrary. But if a loss function is given, then this process of distilling point estimates from the posterior distribution can once more be systematized. For a concrete decision it tells us that parameter value which minimizes the expected loss function under the posterior density function, the so-called “Bayes risk.” This can be considered the Bayesian analog of a point estimate.

For instance, if the loss function is quadratic, then the *posterior mean* is the parameter value which minimizes expected loss.

There is a difference between Bayes risk and the notion of risk we applied previously. The frequentist minimizes expected loss in a large number of repetitions of the trial. This risk is dependent on the unknown parameters, and therefore usually no estimators exist which give minimum risk in all situations. The Bayesian conditions on the data (final criterion!) and minimizes the expected loss where the expectation is taken over the posterior density of the *parameter vector*.

The irreducibility of absence to presences: the absence of knowledge (or a lack of the absence of regularity itself) cannot be represented by a probability distribution. Proof: if I give a certain random variable a neutral prior, then functions of that random variable have non-neutral priors. This argument is made in [Roy97, p. 17]. Many good Bayesians drift away from the subjective point of view and talk about a stratified world: their center of attention is no longer the world out there versus our knowledge of it, but the empirical world versus the underlying systematic forces that shape it.

Bayesians say that frequentists use subjective elements too; their outcomes depend on what the experimenter planned to do, even if he never did it. This argument comes from [Roy97, p. ??]. Nature does not know about the experimenter's plan, and any evidence should be evaluated in a way independent of this.



## CHAPTER 11

## Interval Estimation

Look at our simplest example of an estimator, the sample mean of an independent sample from a normally distributed variable. Since the population mean of a normal variable is at the same time its median, the sample mean will in 50 percent of the cases be larger than the population mean, and in 50 percent of the cases it will be smaller. This is a statement about the procedure how the sample mean was obtained, not about any given observed value of the sample mean. Say in one particular sample the observed sample mean was 3.5. This number 3.5 is either larger or smaller than the true mean, there is no probability involved. But if one were to compute sample means of many different independent samples, then these means would in 50% of the cases lie above and in 50% of the cases below the population mean. This is why one can, from knowing how this one given number was obtained, derive the “confidence” of 50% that the actual mean lies above 3.5, and the same with below. The sample mean can therefore be considered a one-sided confidence bound, although one usually wants higher confidence levels than 50%. (I am 95% confident that  $\phi$  is greater or equal than a certain value computed from the sample.) The concept of “confidence” is nothing but the usual concept of probability if one uses an initial criterion of precision.

The following thought experiment illustrates what is involved. Assume you bought a widget and want to know whether it is defective or not. The obvious way (which would correspond to a “final” criterion of precision) would be to open it up and look if it is defective or not. Now assume we cannot do it: there is no way telling by just looking at it whether it will work. Then another strategy would be to go by an “initial” criterion of precision: we visit the widget factory and look how they make them, how much quality control there is and such. And if we find out that 95% of all widgets coming out of the same factory have no defects, then we have the “confidence” of 95% that our particular widget is not defective either.

The matter becomes only slightly more mystified if one talks about intervals. Again, one should not forget that *confidence intervals are random intervals*. Besides confidence intervals and one-sided confidence bounds one can, if one regards several

parameters simultaneously, also construct confidence rectangles, ellipsoids and more complicated shapes. Therefore we will define in all generality:

Let  $\mathbf{y}$  be a random vector whose distribution depends on some vector of unknown parameters  $\phi \in \Omega$ . A *confidence region* is a prescription which assigns to every possible value  $\mathbf{y}$  of  $\mathbf{y}$  a subset  $R(\mathbf{y}) \subset \Omega$  of parameter space, so that the probability that this subset covers the true value of  $\phi$  is at least a given confidence level  $1 - \alpha$ , i.e.,

$$(11.0.5) \quad \Pr[R(\mathbf{y}) \ni \phi_0 | \phi = \phi_0] \geq 1 - \alpha \quad \text{for all } \phi_0 \in \Omega.$$

The important thing to remember about this definition is that these regions  $R(\mathbf{y})$  are random regions; every time one performs the experiment one obtains a different region.

Now let us go to the specific case of constructing an interval estimate for the parameter  $\mu$  when we have  $n$  independent observations from a normally distributed population  $\sim N(\mu, \sigma^2)$  in which neither  $\mu$  nor  $\sigma^2$  are known. The vector of observations is therefore distributed as  $\mathbf{y} \sim N(\boldsymbol{\nu}\mu, \sigma^2 \mathbf{I})$ , where  $\boldsymbol{\nu}\mu$  is the vector every component of which is  $\mu$ .

I will give you now what I consider to be the cleanest argument deriving the so-called t-interval. It generalizes directly to the F-test in linear regression. It is not the same derivation which you will usually find, and I will bring the usual derivation below for comparison. Recall the observation made earlier, based on (9.1.1), that the sample mean  $\bar{y}$  is that number  $\bar{y} = a$  which minimizes the sum of squared deviations  $\sum (y_i - a)^2$ . (In other words,  $\bar{y}$  is the “least squares estimate” in this situation.) The least squares principle also naturally leads to interval estimates for  $\mu$ : we will say that  $a$  lies in the interval for  $\mu$  if and only if

$$(11.0.6) \quad \frac{\sum (y_i - a)^2}{\sum (y_i - \bar{y})^2} \leq c$$

for some number  $c \geq 1$ . Of course, the value of  $c$  depends on the confidence level, but the beauty of this criterion here is that the value of  $c$  can be determined by the confidence level *alone* without knowledge of the true values of  $\mu$  or  $\sigma^2$ .

To show this, note first that (11.0.6) is equivalent to

$$(11.0.7) \quad \frac{\sum (y_i - a)^2 - \sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \leq c - 1$$

and then apply the identity  $\sum (y_i - a)^2 = \sum (y_i - \bar{y})^2 + n(\bar{y} - a)^2$  to the numerator to get the following equivalent formulation of (11.0.6):

$$(11.0.8) \quad \frac{n(\bar{y} - a)^2}{\sum (y_i - \bar{y})^2} \leq c - 1$$

The confidence level of this interval is the probability that the true  $\mu$  lies in an interval randomly generated using this principle. In other words, it is

$$(11.0.9) \quad \Pr\left[\frac{n(\bar{y} - \mu)^2}{\sum(\mathbf{y}_i - \bar{y})^2} \leq c - 1\right]$$

Although for every *known*  $a$ , the probability that  $a$  lies in the confidence interval depends on the unknown  $\mu$  and  $\sigma^2$ , we will show now that the probability that the *unknown*  $\mu$  lies in the confidence interval does not depend on any unknown parameters. First look at the distribution of the numerator: Since  $\bar{y} \sim N(\mu, \sigma^2/n)$ , it follows  $(\bar{y} - \mu)^2 \sim (\sigma^2/n)\chi_1^2$ . We also know the distribution of the denominator. Earlier we have shown that the variable  $\sum(\mathbf{y}_i - \bar{y})^2$  is a  $\sigma^2\chi_{n-1}^2$ . It is not enough to know the distribution of numerator and denominator separately, we also need their joint distribution. For this go back to our earlier discussion of variance estimation again; there we also showed that  $\bar{y}$  is independent of the vector  $[\mathbf{y}_1 - \bar{y} \ \cdots \ \mathbf{y}_n - \bar{y}]^T$ ; therefore any function of  $\bar{y}$  is also independent of any function of this vector, from which follows that numerator and denominator in our fraction are *independent*. Therefore this fraction is distributed as an  $\sigma^2\chi_1^2$  over an independent  $\sigma^2\chi_{n-1}^2$ , and since the  $\sigma^2$ 's cancel out, this is the same as a  $\chi_1^2$  over an independent  $\chi_{n-1}^2$ . In other words, this distribution does not depend on any unknown parameters!

The definition of a F-distribution with  $k$  and  $m$  degrees of freedom is the distribution of a ratio of a  $\chi_k^2/k$  divided by a  $\chi_m^2/m$ ; therefore if we divide the sum of squares in the numerator by  $n - 1$  we get a F distribution with 1 and  $n - 1$  d.f.:

$$(11.0.10) \quad \frac{(\bar{y} - \mu)^2}{\frac{1}{n-1} \sum(\mathbf{y}_i - \bar{y})^2} \sim F_{1, n-1}$$

If one does not take the square in the numerator, i.e., works with  $\bar{y} - \mu$  instead of  $(\bar{y} - \mu)^2$ , and takes square root in the denominator, one obtains a t-distribution:

$$(11.0.11) \quad \frac{\bar{y} - \mu}{\sqrt{\frac{1}{n} \sum(\mathbf{y}_i - \bar{y})^2}} \sim t_{n-1}$$

The left hand side of this last formula has a suggestive form. It can be written as  $(\bar{y} - \mu)/s_{\bar{y}}$ , where  $s_{\bar{y}}$  is an estimate of the standard deviation of  $\bar{y}$  (it is the square root of the unbiased estimate of the variance of  $\bar{y}$ ). In other words, this t-statistic can be considered an estimate of the number of standard deviations the observed value of  $\bar{y}$  is away from  $\mu$ .

Now we will give, as promised, the usual derivation of the t-confidence intervals, which is based on this interpretation. This usual derivation involves the following two steps:

(1) First assume that  $\sigma^2$  is known. Then it is obvious what to do; for every observation  $\mathbf{y}$  of  $\mathbf{y}$  construct the following interval:

$$(11.0.12) \quad R(\mathbf{y}) = \{u \in \mathbb{R}: |u - \bar{y}| \leq N_{(\alpha/2)}\sigma_{\bar{y}}\}.$$

What do these symbols mean? The interval  $R$  (as in region) has  $\mathbf{y}$  as an argument i.e., it is denoted  $R(\mathbf{y})$ , because it depends on the observed value  $\mathbf{y}$ .  $\mathbb{R}$  is the set of real numbers.  $N_{(\alpha/2)}$  is the upper  $\alpha/2$ -quantile of the Normal distribution, i.e., it is the number  $c$  for which a standard Normal random variable  $z$  satisfies  $\Pr[z \geq c] = \alpha/2$ . Since by the symmetry of the Normal distribution,  $\Pr[z \leq -c] = \alpha/2$  as well, one obtains for a two-sided test:

$$(11.0.13) \quad \Pr[|z| \geq N_{(\alpha/2)}] = \alpha.$$

From this follows the coverage probability:

$$(11.0.14) \quad \Pr[R(\mathbf{y}) \ni \mu] = \Pr[|\mu - \bar{y}| \leq N_{(\alpha/2)}\sigma_{\bar{y}}]$$

$$(11.0.15) \quad = \Pr[|(\mu - \bar{y})/\sigma_{\bar{y}}| \leq N_{(\alpha/2)}] = \Pr[|z| \leq N_{(\alpha/2)}] = 1 - \alpha$$

since  $z = (\bar{y} - \mu)/\sigma_{\bar{y}}$  is a standard Normal. I.e.,  $R(\mathbf{y})$  is a confidence interval for  $\mu$  with confidence level  $1 - \alpha$ .

(2) Second part: what if  $\sigma^2$  is not known? Here a seemingly ad-hoc way would be to replace  $\sigma^2$  by its unbiased estimate  $s^2$ . Of course, then the Normal distribution no longer applies. However if one replaces the normal critical values by those of the  $t_{n-1}$  distribution, one still gets, by miraculous coincidence, a confidence level which is independent of any unknown parameters.

**PROBLEM 183.** If  $\mathbf{y}_i \sim \text{NID}(\mu, \sigma^2)$  (normally independently distributed) with  $n$  observations and  $\sigma^2$  unknown, then the confidence interval for  $\mu$  has the form

$$(11.0.16) \quad R(\mathbf{y}) = \{u \in \mathbb{R}: |u - \bar{y}| \leq t_{(n-1; \alpha/2)} s_{\bar{y}}\}.$$

Here  $t_{(n-1; \alpha/2)}$  is the upper  $\alpha/2$ -quantile of the t distribution with  $n - 1$  degrees of freedom, i.e., it is that number  $c$  for which a random variable  $t$  which has a t distribution with  $n - 1$  degrees of freedom satisfies  $\Pr[t \geq c] = \alpha/2$ . And  $s_{\bar{y}}$  is obtained as follows: write down the standard deviation of  $\bar{y}$  and replace  $\sigma$  by  $s$ . One can also say  $s_{\bar{y}} = \sigma_{\bar{y}} \frac{s}{\sigma}$  where  $\sigma_{\bar{y}}$  is an abbreviated notation for  $\text{std. dev}[\bar{y}] = \sqrt{\text{var}[\bar{y}]}$ .

- a. 1 point Write down the formula for  $s_{\bar{y}}$ .

ANSWER. Start with  $\sigma_{\bar{y}}^2 = \text{var}[\bar{y}] = \frac{\sigma^2}{n}$ , therefore  $\sigma_{\bar{y}} = \sigma/\sqrt{n}$ , and

$$(11.0.17) \quad s_{\bar{y}} = s/\sqrt{n} = \sqrt{\frac{\sum(\mathbf{y}_i - \bar{y})^2}{n(n-1)}}$$

- b. 2 points Compute the coverage probability of the interval (11.0.16).

TABLE 1. Percentiles of Student's t Distribution. Table entry  $x$  satisfies  $\Pr[t_n \leq x] = p$ .

$n$	$p =$					
	.750	.900	.950	.975	.990	.995
1	1.000	3.078	6.314	12.706	31.821	63.657
2	0.817	1.886	2.920	4.303	6.965	9.925
3	0.765	1.638	2.354	3.182	4.541	5.841
4	0.741	1.533	2.132	2.776	3.747	4.604
5	0.727	1.476	2.015	2.571	3.365	4.032

ANSWER. The coverage probability is

$$(11.0.18) \quad \Pr[\mathbf{R}(\mathbf{y}) \ni \mu] = \Pr\left[|\mu - \bar{y}| \leq t_{(n-1; \alpha/2)} s_{\bar{y}}\right]$$

$$(11.0.19) \quad = \Pr\left[\left|\frac{\mu - \bar{y}}{s_{\bar{y}}}\right| \leq t_{(n-1; \alpha/2)}\right]$$

$$(11.0.20) \quad = \Pr\left[\left|\frac{(\mu - \bar{y})/\sigma_{\bar{y}}}{s_{\bar{y}}/\sigma_{\bar{y}}}\right| \leq t_{(n-1; \alpha/2)}\right]$$

$$(11.0.21) \quad = \Pr\left[\left|\frac{(\bar{y} - \mu)/\sigma_{\bar{y}}}{s/\sigma}\right| \leq t_{(n-1; \alpha/2)}\right]$$

$$(11.0.22) \quad = 1 - \alpha,$$

because the expression in the numerator is a standard normal, and the expression in the denominator is the square root of an independent  $\chi_{n-1}^2$  divided by  $n - 1$ . The random variable between the absolute signs has therefore a t-distribution, and (11.0.22) follows from (30.4.8).  $\square$

• c. 2 points Four independent observations are available of a normal random variable with unknown mean  $\mu$  and variance  $\sigma^2$ : the values are  $-2, -\sqrt{2}, +\sqrt{2},$  and  $+2$ . (These are not the kind of numbers you are usually reading off a measurement instrument, but they make the calculation easy). Give a 95% confidence interval for  $\mu$ . Table 1 gives the percentiles of the t-distribution.

ANSWER. In our situation

$$(11.0.23) \quad \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_3$$

According to table 1, for  $b = 3.182$  follows

$$(11.0.24) \quad \Pr[t_3 \leq b] = 0.975$$

therefore

$$(11.0.25) \quad \Pr[t_3 > b] = 0.025$$

and by symmetry of the t-distribution

$$(11.0.26) \quad \Pr[t_3 < -b] = 0.025$$

Now subtract (11.0.26) from (11.0.24) to get

$$(11.0.27) \quad \Pr[-b \leq t_3 \leq b] = 0.95$$

or

$$(11.0.28) \quad \Pr[|t_3| \leq b] = 0.95$$

or, plugging in the formula for  $t_3$ ,

$$(11.0.29) \quad \Pr\left[\left|\frac{\bar{x} - \mu}{s/\sqrt{n}}\right| \leq b\right] = .95$$

$$(11.0.30) \quad \Pr[|\bar{x} - \mu| \leq bs/\sqrt{n}] = .95$$

$$(11.0.31) \quad \Pr[-bs/\sqrt{n} \leq \mu - \bar{x} \leq bs/\sqrt{n}] = .95$$

$$(11.0.32) \quad \Pr[\bar{x} - bs/\sqrt{n} \leq \mu \leq \bar{x} + bs/\sqrt{n}] = .95$$

the confidence interval is therefore  $[\bar{x} - bs/\sqrt{n}, \bar{x} + bs/\sqrt{n}]$ . In our sample,  $\bar{x} = 0, s^2 = \frac{12}{3} = 4, n = 4,$  therefore  $s^2/n = 1,$  therefore also  $s/\sqrt{n} = 1$ . So the sample value of the confidence interval is  $[-3.182, +3.182]$ .

PROBLEM 184. Using R, construct 20 samples of 12 observation each from  $N(0, 1)$  distribution, construct the 95% confidence t-intervals for the mean based on these 20 samples, plot these intervals, and count how many intervals contain the true mean.

Here are the commands: `stdnorms<-matrix(rnorm(240),nrow=12,ncol=20)` a  $12 \times 20$  matrix containing 240 independent random normals. You get the vector containing the midpoints of the confidence intervals by the assignment `midpts=apply(stdnorms,2,mean)`. About `apply` see [BCW96, p. 130]. The vector containing the half width of each confidence interval can be obtained by another use of `apply`: `halfwidth <- (qt(0.975,11)/sqrt(12)) * sqrt(apply(stdnorms,2,var))`; print the values on the screen you may simply issue the command `cbind(midpts-halfwidth, midpts+halfwidth)`. But it is much better to plot them. Since such a plot does not have one of the usual formats, we have to put it together with some low-level commands. See [BCW96, page 325]. At the very minimum we need the following: `frame()` starts a new plot; `par(usr = c(1,20, range(c(midpts-halfwidth, midpts+halfwidth))))` sets a coordinate system which accommodates all intervals. The 20 confidence intervals are plotted by `segments(1:20, midpts-halfwidth, 1:20, midpts+halfwidth)`. Finally, `abline(0,0)` adds a horizontal line, so that you can see how many intervals contain the true mean.

The `ecmet` package has a function `confint.segments` which draws such plots automatically. Choose how many observations in each experiment (the argument

$n$ ), and how many confidence intervals (the argument `rep`), and the confidence level `level` (the default is here 95%), and then issue, e.g. the command `confint.segments(n=`

Here is the transcript of the function:

```
confint.segments <- function(n, rep, level = 95/100)
{
  stdnormals <- matrix(rnorm(n * rep), nrow = n, ncol = rep)
  midpts <- apply(stdnormals, 2, mean)
  halfwidth <- qt(p=(1 + level)/2, df= n - 1) * sqrt(1/n)* sqrt(apply(stdnormals, 2,
  frame()
  x <- c(1:rep, 1:rep)
  y <- c(midpts + halfwidth, midpts - halfwidth)
  par(usr = c(1, rep, range(y)))
  segments(1:rep, midpts - halfwidth, 1:rep, midpts + halfwidth)
  abline(0, 0)
  invisible(cbind(x,y))
}
```

This function draws the plot as a “side effect,” but it also returns a matrix with the coordinates of the endpoints of the plots (without printing them on the screen). This matrix can be used as input for the `identify` function. If you do for instance `iddata<-confint.segments(12,20)` and then `identify(iddata, labels=iddata[,2])`, then the following happens: if you move the mouse cursor on the graph near one of the endpoints of one of the intervals, and click the left button, then it will print on the graph the coordinate of the boundary of this interval. Clicking any other button of the mouse gets you out of the `identify` function.

## CHAPTER 12

## Hypothesis Testing

Imagine you are a business person considering a major investment in order to launch a new product. The sales prospects of this product are not known with certainty. You have to rely on the outcome of  $n$  marketing surveys that measure the demand for the product once it is offered. If  $\mu$  is the actual (unknown) rate of return on the investment, each of these surveys here will be modeled as a random variable, which has a Normal distribution with this mean  $\mu$  and known variance 1. Let  $y_1, y_2, \dots, y_n$  be the observed survey results. How would you decide whether to build the plant?

The intuitively reasonable thing to do is to go ahead with the investment if the sample mean of the observations is greater than a given value  $c$ , and not to do it otherwise. This is indeed an optimal decision rule, and we will discuss in what respect it is, and how  $c$  should be picked.

Your decision can be the wrong decision in two different ways: either you decide to go ahead with the investment although there will be no demand for the product, or you fail to invest although there would have been demand. There is no decision rule which eliminates both errors at once; the first error would be minimized by the rule never to produce, and the second by the rule always to produce. In order to determine the right tradeoff between these errors, it is important to be aware of their *asymmetry*. The error to go ahead with production although there is no demand has potentially disastrous consequences (loss of a lot of money), while the other error may cause you to miss a profit opportunity, but there is no actual loss involved, and presumably you can find other opportunities to invest your money.

To express this asymmetry, the error with the potentially disastrous consequences is called “error of type one,” and the other “error of type two.” The distinction between type one and type two errors can also be made in other cases. Locking up an innocent person is an error of type one, while letting a criminal go unpunished is an error of type two; publishing a paper with false results is an error of type one, while foregoing an opportunity to publish is an error of type two (at least this is what it ought to be).

Such an asymmetric situation calls for an asymmetric decision rule. One needs strict safeguards against committing an error of type one, and if there are several decision rules which are equally safe with respect to errors of type one, then one would select among those that decision rule which minimizes the error of type two.

Let us look here at decision rules of the form: make the investment if  $\bar{y} > c$ . An error of type one occurs if the decision rule advises you to make the investment while there is no demand for the product. This will be the case if  $\bar{y} > c$  but  $\mu \leq 0$ . The probability of this error depends on the unknown parameter  $\mu$ , but it is at most  $\alpha = \Pr[\bar{y} > c | \mu = 0]$ . This maximum value of the type one error probability is called the significance level, and you, as the director of the firm, will have to decide on it depending on how tolerable it is to lose money on this venture, which presumably depends on the chances to lose money on alternative investments. It is a serious shortcoming of the classical theory of hypothesis testing that it does not provide good guidelines how  $\alpha$  should be chosen, and how it should change with sample size. Instead, there is the tradition to choose  $\alpha$  to be either 5% or 1% or 0.1%. Given a table of the cumulative standard normal distribution function allows you to find that  $c$  for which  $\Pr[\bar{y} > c | \mu = 0] = \alpha$ .

**PROBLEM 185.** *2 points* Assume each  $y_i \sim N(\mu, 1)$ ,  $n = 400$  and  $\alpha = 0.05$ , and different  $y_i$  are independent. Compute the value  $c$  which satisfies  $\Pr[\bar{y} > c | \mu = 0] = \alpha$ . You should either look it up in a table and include a xerox copy of the table with the entry circled and the complete bibliographic reference written on the xerox copy, or do it on a computer, writing exactly which commands you used. In R, the function `qnorm` does what you need, find out about it by typing `help(qnorm)`.

**ANSWER.** In the case  $n = 400$ ,  $\bar{y}$  has variance  $1/400$  and therefore standard deviation  $1/20$ .  $20\bar{y}$  is a standard normal: from  $\Pr[\bar{y} > c | \mu = 0] = 0.05$  follows  $\Pr[20\bar{y} > 20c | \mu = 0] = 0.05$ . Therefore  $20c = 1.645$  can be looked up in a table, perhaps use [JHG<sup>+</sup>88, p. 986], row for  $\infty$  d.f.

Let us do this in R. The  $p$ -“quantile” of the distribution of the random variable  $y$  is defined as that value  $q$  for which  $\Pr[y \leq q] = p$ . If  $y$  is normally distributed, this quantile is computed by the R-function `qnorm(p, mean=0, sd=1, lower.tail=TRUE)`. In the present case we need either `qnorm(p=1-0.05, mean=0, sd=0.05)` or `qnorm(p=0.05, mean=0, sd=0.05, lower.tail=FALSE)` which gives the value 0.08224268.

Choosing a decision which makes a loss unlikely is not enough; your decision must also give you a chance of success. E.g., the decision rule to build the plant if  $-0.06 \leq \bar{y} \leq -0.05$  and not to build it otherwise is completely perverse, although the significance level of this decision rule is approximately 4% (if  $n = 100$ ). In other words, the significance level is not enough information for evaluating the performance of the test. You also need the “power function,” which gives you the probability with which the test advises you to make the “critical” decision, as a function

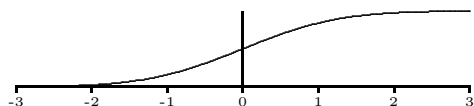


FIGURE 1. Eventually this Figure will show the Power function of a one-sided normal test, i.e., the probability of error of type one as a function of  $\mu$ ; right now this is simply the cdf of a Standard Normal

the true parameter values. (Here the “critical” decision is that decision which might potentially lead to an error of type one.) By the definition of the significance level, the power function does not exceed the significance level for those parameter values for which going ahead would lead to a type 1 error. But only those tests are “powerful” whose power function is high for those parameter values for which it would be correct to go ahead. In our case, the power function must be below 0.05 when  $\mu \leq 0$ , and we want it as high as possible when  $\mu > 0$ . Figure 1 shows the power function for the decision rule to go ahead whenever  $\bar{y} \geq c$ , where  $c$  is chosen in such a way that the significance level is 5%, for  $n = 100$ .

The hypothesis whose *rejection*, although it is true, constitutes an error of type one, is called the *null hypothesis*, and its alternative the *alternative hypothesis*. (In the examples the null hypotheses were: the return on the investment is zero or negative, the defendant is innocent, or the results about which one wants to publish a research paper are wrong.) The null hypothesis is therefore the hypothesis that nothing is the case. The test tests whether this hypothesis should be rejected, will safeguard against the hypothesis one *wants to reject* but one is afraid to reject *erroneously*. If you reject the null hypothesis, you don’t want to regret it.

Mathematically, every test can be identified with its null hypothesis, which is a region in parameter space (often consisting of one point only), and its “critical region,” which is the *event* that the test comes out in favor of the “critical decision,” i.e., rejects the null hypothesis. The critical region is usually an event of the form that the value of a certain random variable, the “test statistic,” is within a given range, usually that it is too high. The power function of the test is the probability of the critical region as a function of the unknown parameters, and the significance level is the maximum (or, if this maximum depends on unknown parameters, any upper bound) of the power function over the null hypothesis.

**PROBLEM 186.** *Mr. Jones is on trial for counterfeiting Picasso paintings, and you are an expert witness who has developed fool-proof statistical significance tests for identifying the painter of a given painting.*

- a. 2 points *There are two ways you can set up your test.*

- a: *You can either say: The null hypothesis is that the painting was done by Picasso, and the alternative hypothesis that it was done by Mr. Jones.*
- b: *Alternatively, you might say: The null hypothesis is that the painting was done by Mr. Jones, and the alternative hypothesis that it was done by Picasso.*

*Does it matter which way you do the test, and if so, which way is the correct one? Give a reason to your answer, i.e., say what would be the consequences of testing the incorrect way.*

**ANSWER.** The determination of what the null and what the alternative hypothesis is depends on what is considered to be the catastrophic error which is to be guarded against. On a trial, Mr. Jones is considered innocent until proven guilty. Mr. Jones should not be convicted unless he can be proven guilty beyond “reasonable doubt.” Therefore the test must be set up in such a way that the hypothesis that the painting is by Picasso will only be rejected if the chance that it is actually by Picasso is very small. The error of type one is that the painting is considered counterfeited although it is really by Picasso. Since the error of type one is always the error to reject the null hypothesis although it is true, solution a. is the correct one. You are not proving, you are testing.

- b. 2 points *After the trial a customer calls you who is in the process of acquiring a very expensive alleged Picasso painting, and who wants to be sure that this painting is not one of Jones’s falsifications. Would you now set up your test in the same way as in the trial or in the opposite way?*

**ANSWER.** It is worse to spend money on a counterfeit painting than to forego purchasing a true Picasso. Therefore the null hypothesis would be that the painting was done by Mr. Jones, and the alternative is the opposite way.

### 12.1. Duality between Significance Tests and Confidence Regions

There is a duality between confidence regions with confidence level  $1 - \alpha$  and certain significance tests. Let us look at a family of significance tests, which all have a significance level  $\leq \alpha$ , and which define for every possible value of the parameter  $\phi_0 \in \Omega$  a critical region  $C(\phi_0)$  for rejecting the simple null hypothesis that the true parameter is equal to  $\phi_0$ . The condition that all significance levels are  $\leq \alpha$  means mathematically

$$(12.1.1) \quad \Pr[C(\phi_0) | \phi = \phi_0] \leq \alpha \quad \text{for all } \phi_0 \in \Omega.$$

Mathematically, confidence regions and such families of tests are one and the same thing: if one has a confidence region  $R(\mathbf{y})$ , one can define a test of the null hypothesis  $\phi = \phi_0$  as follows: for an observed outcome  $\mathbf{y}$  reject the null hypothesis if and only if  $\phi_0$  is not contained in  $R(\mathbf{y})$ . On the other hand, given a family of tests one can build a confidence region by the prescription:  $R(\mathbf{y})$  is the set of all the parameter values which would not be rejected by a test based on observation  $\mathbf{y}$ .

PROBLEM 187. Show that with these definitions, equations (11.0.5) and (12.1.1) are equivalent.

ANSWER. Since  $\phi_0 \in R(\mathbf{y})$  iff  $\mathbf{y} \in C'(\phi_0)$  (the complement of the critical region rejecting that the parameter value is  $\phi_0$ ), it follows  $\Pr[R(\mathbf{y}) \in \phi_0 | \phi = \phi_0] = 1 - \Pr[C(\phi_0) | \phi = \phi_0] \geq 1 - \alpha$ .  $\square$

This duality is discussed in [BD77, pp. 177–182].

### 12.2. The Neyman Pearson Lemma and Likelihood Ratio Tests

Look one more time at the example with the fertilizer. Why are we considering only regions of the form  $\bar{y} \geq \mu_0$ , why not one of the form  $\mu_1 \leq \bar{y} \leq \mu_2$ , or maybe not use the mean but decide to build if  $\mathbf{y}_1 \geq \mu_3$ ? Here the  $\mu_1, \mu_2$ , and  $\mu_3$  can be chosen such that the probability of committing an error of type one is still  $\alpha$ .

It seems intuitively clear that these alternative decision rules are not reasonable. The Neyman Pearson lemma proves this intuition right. It says that the critical regions of the form  $\bar{y} \geq \mu_0$  are uniformly most powerful, in the sense that every other critical region with same probability of type one error has equal or higher probability of committing error of type two, regardless of the true value of  $\mu$ .

Here are formulation and proof of the Neyman Pearson lemma, first for the case that both null hypothesis and alternative hypothesis are simple:  $H_0 : \theta = \theta_0, H_A : \theta = \theta_1$ . In other words, we want to determine on the basis of the observations of the random variables  $\mathbf{y}_1, \dots, \mathbf{y}_n$  whether the true  $\theta$  was  $\theta_0$  or  $\theta_1$ , and a determination  $\theta = \theta_1$  when in fact  $\theta = \theta_0$  is an error of type one. The critical region C is the set of all outcomes that lead us to conclude that the parameter has value  $\theta_1$ .

The Neyman Pearson lemma says that a uniformly most powerful test exists in this situation. It is a so-called likelihood-ratio test, which has the following critical region:

$$(12.2.1) \quad C = \{y_1, \dots, y_n : L(y_1, \dots, y_n; \theta_1) \geq kL(y_1, \dots, y_n; \theta_0)\}.$$

C consists of those outcomes for which  $\theta_1$  is at least  $k$  times as likely as  $\theta_0$  (where  $k$  is chosen such that  $\Pr[C|\theta_0] = \alpha$ ).

To prove that this decision rule is uniformly most powerful, assume D is the critical region of a different test with same significance level  $\alpha$ , i.e., if the null hypothesis is correct, then C and D reject (and therefore commit an error of type one) with equally low probabilities  $\alpha$ . In formulas,  $\Pr[C|\theta_0] = \Pr[D|\theta_0] = \alpha$ . Look at figure 2 with  $C = U \cup V$  and  $D = V \cup W$ . Since C and D have the same significance level, it follows

$$(12.2.2) \quad \Pr[U|\theta_0] = \Pr[W|\theta_0].$$

FIGURE 2. Venn Diagram for Proof of Neyman Pearson Lemma ec660.1005

Also

$$(12.2.3) \quad \Pr[U|\theta_1] \geq k \Pr[U|\theta_0],$$

since  $U \subset C$  and C were chosen such that the likelihood (density) function of the alternative hypothesis is high relatively to that of the null hypothesis. Since  $W$  is outside C, the same argument gives

$$(12.2.4) \quad \Pr[W|\theta_1] \leq k \Pr[W|\theta_0].$$

Linking those two inequalities and the equality gives

$$(12.2.5) \quad \Pr[W|\theta_1] \leq k \Pr[W|\theta_0] = k \Pr[U|\theta_0] \leq \Pr[U|\theta_1],$$

hence  $\Pr[D|\theta_1] \leq \Pr[C|\theta_1]$ . In other words, if  $\theta_1$  is the correct parameter value, then C will discover this and reject at least as often as D. Therefore C is at least as powerful as D, or the type two error probability of C is at least as small as that of D.

Back to our fertilizer example. To make both null and alternative hypothesis simple, assume that either  $\mu = 0$  (fertilizer is ineffective) or  $\mu = t$  for some fixed

$t > 0$ . Then the likelihood ratio critical region has the form

$$(12.2.6) \quad C = \{y_1, \dots, y_n : \left(\frac{1}{\sqrt{2\pi}}\right)^n e^{-\frac{1}{2}((y_1-t)^2 + \dots + (y_n-t)^2)} \geq k \left(\frac{1}{\sqrt{2\pi}}\right)^n e^{-\frac{1}{2}(y_1^2 + \dots + y_n^2)}\}$$

$$(12.2.7) \quad = \{y_1, \dots, y_n : -\frac{1}{2}((y_1-t)^2 + \dots + (y_n-t)^2) \geq \ln k - \frac{1}{2}(y_1^2 + \dots + y_n^2)\}$$

$$(12.2.8) \quad = \{y_1, \dots, y_n : t(y_1 + \dots + y_n) - \frac{t^2 n}{2} \geq \ln k\}$$

$$(12.2.9) \quad = \{y_1, \dots, y_n : \bar{y} \geq \frac{\ln k}{nt} + \frac{t}{2}\}$$

i.e.,  $C$  has the form  $\bar{y} \geq$  some constant. The dependence of this constant on  $k$  is not relevant, since this constant is usually chosen such that the maximum probability of error of type one is equal to the given significance level.

PROBLEM 188. *8 points* You have four independent observations  $y_1, \dots, y_4$  from an  $N(\mu, 1)$ , and you are testing the null hypothesis  $\mu = 0$  against the alternative hypothesis  $\mu = 1$ . For your test you are using the likelihood ratio test with critical region

$$(12.2.10) \quad C = \{y_1, \dots, y_4 : L(y_1, \dots, y_4; \mu = 1) \geq 3.633 \cdot L(y_1, \dots, y_4; \mu = 0)\}.$$

Compute the significance level of this test. (According to the Neyman-Pearson lemma, this is the uniformly most powerful test for this significance level.) Hints: In order to show this you need to know that  $\ln 3.633 = 1.29$ , everything else can be done without a calculator. Along the way you may want to show that  $C$  can also be written in the form  $C = \{y_1, \dots, y_4 : y_1 + \dots + y_4 \geq 3.290\}$ .

ANSWER. Here is the equation which determines when  $y_1, \dots, y_4$  lie in  $C$ :

$$(12.2.11) \quad (2\pi)^{-2} \exp -\frac{1}{2}((y_1-1)^2 + \dots + (y_4-1)^2) \geq 3.633 \cdot (2\pi)^{-2} \exp -\frac{1}{2}(y_1^2 + \dots + y_4^2)$$

$$(12.2.12) \quad -\frac{1}{2}((y_1-1)^2 + \dots + (y_4-1)^2) \geq \ln(3.633) - \frac{1}{2}(y_1^2 + \dots + y_4^2)$$

$$(12.2.13) \quad y_1 + \dots + y_4 - 2 \geq 1.290$$

Since  $\Pr\{y_1 + \dots + y_4 \geq 3.290\} = \Pr\{z = (y_1 + \dots + y_4)/2 \geq 1.645\}$  and  $z$  is a standard normal, one obtains the significance level of 5% from the standard normal table or the t-table.  $\square$

Note that due to the properties of the Normal distribution, this critical region, for a given significance level, does not depend at all on the value of  $t$ . Therefore this test is uniformly most powerful against the composite hypothesis  $\mu > 0$ .

One can also write the null hypothesis as the composite hypothesis  $\mu \leq 0$ , because the highest probability of type one error will still be attained when  $\mu = 0$ . This completes the proof that the test given in the original fertilizer example is uniformly most powerful.

Most other distributions discussed here are equally well behaved, therefore uniformly most powerful one-sided tests exist not only for the mean of a normal with known variance, but also the variance of a normal with known mean, or the parameters of a Bernoulli and Poisson distribution.

However the given one-sided hypothesis is the only situation in which a uniformly most powerful test exists. In other situations, the *generalized likelihood ratio test* has good properties even though it is no longer uniformly most powerful. Many known tests (e.g., the F test) are generalized likelihood ratio tests.

Assume you want to test the composite null hypothesis  $H_0 : \theta \in \omega$ , where  $\omega$  is a subset of the parameter space, against the alternative  $H_A : \theta \in \Omega$ , where  $\Omega \supset \omega$  is a more comprehensive subset of the parameter space.  $\omega$  and  $\Omega$  are defined by functions with continuous first-order derivatives. The generalized likelihood ratio critical region has the form

$$(12.2.14) \quad C = \{x_1, \dots, x_n : \frac{\sup_{\theta \in \Omega} L(x_1, \dots, x_n; \theta)}{\sup_{\theta \in \omega} L(x_1, \dots, x_n; \theta)} \geq k\}$$

where  $k$  is chosen such that the probability of the critical region when the null hypothesis is true has as its maximum the desired significance level. It can be shown that twice the log of this quotient is asymptotically distributed as a  $\chi^2_{q-s}$ , where  $q$  is the dimension of  $\Omega$  and  $s$  the dimension of  $\omega$ . (Sometimes the likelihood ratio is defined as the inverse of this ratio, but whenever possible we will define our test statistics so that the null hypothesis is rejected if the value of the test statistic is too large.)

In order to perform a likelihood ratio test, the following steps are necessary. First construct the MLE's for  $\theta \in \Omega$  and  $\theta \in \omega$ , then take twice the difference of the attained levels of the log likelihood functions, and compare with the  $\chi^2$  tables.

### 12.3. The Wald, Likelihood Ratio, and Lagrange Multiplier Tests

Let us start with the generalized Wald test. Assume  $\tilde{\theta}$  is an asymptotically normal estimator of  $\theta$ , whose asymptotic distribution is  $N(\theta, \Psi)$ . Assume furthermore that  $\hat{\Psi}$  is a consistent estimate of  $\Psi$ . Then the following statistic is called the generalized Wald statistic. It can be used for an asymptotic test of the hypothesis  $\mathbf{h}(\theta) = \mathbf{o}$ , where  $\mathbf{h}$  is a  $q$ -vector-valued differentiable function:

$$(12.3.1) \quad \text{G.W.} = \mathbf{h}(\tilde{\theta})^\top \left\{ \frac{\partial \mathbf{h}}{\partial \theta^\top} \Big|_{\tilde{\theta}} \hat{\Psi} \frac{\partial \mathbf{h}}{\partial \theta} \Big|_{\tilde{\theta}} \right\}^{-1} \mathbf{h}(\tilde{\theta})$$



Under the null hypothesis, this test statistic is asymptotically distributed as a  $\chi_q^2$ . To understand this, note that for all  $\boldsymbol{\theta}$  close to  $\tilde{\boldsymbol{\theta}}$ ,  $\mathbf{h}(\boldsymbol{\theta}) \asymp \mathbf{h}(\tilde{\boldsymbol{\theta}}) + \frac{\partial \mathbf{h}}{\partial \boldsymbol{\theta}^\top} \Big|_{\tilde{\boldsymbol{\theta}}} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})$ . Taking covariances

$$(12.3.2) \quad \frac{\partial \mathbf{h}}{\partial \boldsymbol{\theta}^\top} \Big|_{\hat{\boldsymbol{\theta}}} \hat{\boldsymbol{\Psi}} \frac{\partial \mathbf{h}^\top}{\partial \boldsymbol{\theta}} \Big|_{\hat{\boldsymbol{\theta}}}$$

is an estimate of the covariance matrix of  $\mathbf{h}(\tilde{\boldsymbol{\theta}})$ . I.e., one takes  $\mathbf{h}(\tilde{\boldsymbol{\theta}})$  twice and “divides” it by its covariance matrix.

Now let us make more stringent assumptions. Assume the density  $f_{\mathbf{x}}(\mathbf{x}; \boldsymbol{\theta})$  of  $\mathbf{x}$  depends on the parameter vector  $\boldsymbol{\theta}$ . We are assuming that the conditions are satisfied which ensure asymptotic normality of the maximum likelihood estimator  $\hat{\boldsymbol{\theta}}$  and also of  $\bar{\boldsymbol{\theta}}$ , the constrained maximum likelihood estimator subject to the constraint  $\mathbf{h}(\boldsymbol{\theta}) = \mathbf{o}$ .

There are three famous tests to test this hypothesis, which asymptotically are all distributed like  $\chi_q^2$ . The likelihood-ratio test is

$$(12.3.3) \quad LRT = -2 \log \frac{\max_{\mathbf{h}(\boldsymbol{\theta})=\mathbf{o}} f_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\theta})}{\max_{\boldsymbol{\theta}} f_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\theta})} = 2(\log f_{\mathbf{y}}(\mathbf{y}, \hat{\boldsymbol{\theta}}) - \log f_{\mathbf{y}}(\mathbf{y}, \bar{\boldsymbol{\theta}}))$$

It rejects if imposing the constraint reduces the attained level of the likelihood function too much.

The Wald test has the form

$$(12.3.4) \quad \text{Wald} = -\mathbf{h}(\hat{\boldsymbol{\theta}})^\top \left\{ \frac{\partial \mathbf{h}}{\partial \boldsymbol{\theta}^\top} \Big|_{\hat{\boldsymbol{\theta}}} \left( \frac{\partial^2 \log f(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\hat{\boldsymbol{\theta}}} \right)^{-1} \frac{\partial \mathbf{h}^\top}{\partial \boldsymbol{\theta}} \Big|_{\hat{\boldsymbol{\theta}}} \right\}^{-1} \mathbf{h}(\hat{\boldsymbol{\theta}})$$

To understand this formula, note that  $-\left( \mathcal{E} \left[ \frac{\partial^2 \log f(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right] \right)^{-1}$  is the Cramer Rao lower bound, and since all maximum likelihood estimators asymptotically attain the CRLB, it is the asymptotic covariance matrix of  $\hat{\boldsymbol{\theta}}$ . If one does not take the expected value but plugs  $\hat{\boldsymbol{\theta}}$  into these partial derivatives of the log likelihood function, one gets a consistent estimate of the asymptotic covariance matrix. Therefore the Wald test is a special case of the generalized Wald test.

Finally the score test has the form

$$(12.3.5) \quad \text{Score} = -\frac{\partial \log f(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\bar{\boldsymbol{\theta}}} \left( \frac{\partial^2 \log f(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\bar{\boldsymbol{\theta}}} \right)^{-1} \frac{\partial \log f(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\bar{\boldsymbol{\theta}}}$$

This test tests whether the score, i.e., the gradient of the unconstrained log likelihood function, evaluated at the constrained maximum likelihood estimator, is too far away from zero. To understand this formula, remember that we showed in the proof of the Cramer-Rao lower bound that the negative of the expected value of the Hessian

$-\mathcal{E} \left[ \frac{\partial^2 \log f(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right]$  is the covariance matrix of the score, i.e., here we take the score twice and divide it by its estimated covariance matrix.

## CHAPTER 13

## General Principles of Econometric Modelling

[Gre97, 6.1 on p. 220] says: “An econometric study begins with a set of propositions about some aspect of the economy. The theory specifies a set of precise, deterministic relationships among variables. Familiar examples are demand equations, production functions, and macroeconomic models. The empirical investigation provides estimates of unknown parameters in the model, such as elasticities or the marginal propensity to consume, and usually attempts to measure the validity of the theory against the behavior of the observable data.”

[Hen95, p. 6] distinguishes between two extremes: “‘Theory-driven’ approaches, in which the model is derived from a priori theory and calibrated from data evidence. They suffer from theory dependence in that their credibility depends on the credibility of the theory from which they arose—when that theory is discarded, so is the associated evidence.” The other extreme is “‘Data-driven’ approaches, where models are developed to closely describe the data . . . These suffer from sample dependence in that accidental and transient data features are embodied as tightly in the model as permanent aspects, so that extension of the data set often reveal predictive failure.”

Hendry proposes the following useful distinction of 4 levels of knowledge:

*A* Consider the situation where we know the complete structure of the process which generates economic data and the values of all its parameters. This is the equivalent of a probability theory course (example: rolling a perfect die), but involves economic theory and econometric concepts.

*B* consider a known economic structure with unknown values of the parameters. Equivalent to an estimation and inference course in statistics (example: independent rolls of an imperfect die and estimating the probabilities of the different faces) but focusing on econometrically relevant aspects.

*C* is “the empirically relevant situation where neither the form of the data-generating process nor its parameter values are known. (Here one does not know whether the rolls of the die are independent, or whether the probabilities of the different faces remain constant.) Model discovery, evaluation, data mining, model-search procedures, and associated methodological issues.

*D* Forecasting the future when the data outcomes are unknown. (Model of money demand under financial innovation).

The example of Keynes’s consumption function in [Gre97, pp. 221/22] sounds like the beginning as if it was close to *B*, but in the further discussion Greene goes more and more over to *C*. It is remarkable here that economic theory usually does not yield functional forms. Greene then says: the most common functional form is the linear one  $c = \alpha + \beta x$  with  $\alpha > 0$  and  $0 < \beta < 1$ . He does not mention the aggregation problem hidden in this. Then he says: “But the linear function is only approximate. In fact, it is unlikely that consumption and income can be connected by any simple relationship. The **deterministic** relationship is clearly inadequate.” Here Greene uses a random relationship to model a relationship which is quantitatively “fuzzy.” This is an interesting and relevant application of randomness.

A sentence later Green backtracks from this insight and says: “We are not so ambitious as to attempt to capture every influence in the relationship, but only those that are substantial enough to model directly.” The “fuzziness” is not due to a lack of ambition of the researcher, but the world is inherently quantitatively fuzzy. It is not that we don’t know the law, but there is no law; not everything that happens in an economy is driven by economic laws. Greene’s own example, in Figure 6.2, that during the war years consumption was below the trend line, shows this.

Greene’s next example is the relationship between income and education. This illustrates multiple instead of simple regression: one must also include age, and the square of age, even if one is not interested in the effect which age has, but in order to “control” for this effect, so that the effects of education and age will not be confounded.

**PROBLEM 189.** *Why should a regression of income on education include not only age but also the square of age?*

**ANSWER.** Because the effect of age becomes smaller with increases in age.

Critical Realist approaches are [Ron02] and [Mor02].

or

$$(14.1.3) \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad \boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2 \mathbf{I})$$

PROBLEM 190. 1 point Compute the matrix product

$$\begin{bmatrix} 1 & 2 & 5 \\ 0 & 3 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 2 & 1 \\ 3 & 8 \end{bmatrix}$$

ANSWER.

$$\begin{bmatrix} 1 & 2 & 5 \\ 0 & 3 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 2 & 1 \\ 3 & 8 \end{bmatrix} = \begin{bmatrix} 1 \cdot 4 + 2 \cdot 2 + 5 \cdot 3 & 1 \cdot 0 + 2 \cdot 1 + 5 \cdot 8 \\ 0 \cdot 4 + 3 \cdot 2 + 1 \cdot 3 & 0 \cdot 0 + 3 \cdot 1 + 1 \cdot 8 \end{bmatrix} = \begin{bmatrix} 23 & 42 \\ 9 & 11 \end{bmatrix}$$

## CHAPTER 14

## Mean-Variance Analysis in the Linear Model

In the present chapter, the only distributional assumptions are that means and variances exist. (From this follows that also the covariances exist).

## 14.1. Three Versions of the Linear Model

As background reading please read [CD97, Chapter 1].

Following [JHG<sup>+</sup>88, Chapter 5], we will start with three different linear statistical models. Model 1 is the simplest estimation problem already familiar from chapter 9, with  $n$  independent observations from the same distribution, call them  $\mathbf{y}_1, \dots, \mathbf{y}_n$ . The only thing known about the distribution is that mean and variance exist, call them  $\boldsymbol{\mu}$  and  $\sigma^2$ . In order to write this as a special case of the “linear model,” define  $\boldsymbol{\varepsilon}_i = \mathbf{y}_i - \boldsymbol{\mu}$ , and define the vectors  $\mathbf{y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_n]^\top$ ,  $\boldsymbol{\varepsilon} = [\boldsymbol{\varepsilon}_1 \ \boldsymbol{\varepsilon}_2 \ \dots \ \boldsymbol{\varepsilon}_n]^\top$ , and  $\boldsymbol{\iota} = [1 \ 1 \ \dots \ 1]^\top$ . Then one can write the model in the form

$$(14.1.1) \quad \mathbf{y} = \boldsymbol{\iota}\boldsymbol{\mu} + \boldsymbol{\varepsilon} \quad \boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2 \mathbf{I})$$

The notation  $\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2 \mathbf{I})$  is shorthand for  $\mathcal{E}[\boldsymbol{\varepsilon}] = \mathbf{o}$  (the null vector) and  $\mathcal{V}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}$  ( $\sigma^2$  times the identity matrix, which has 1’s in the diagonal and 0’s elsewhere).  $\boldsymbol{\mu}$  is the deterministic part of all the  $\mathbf{y}_i$ , and  $\boldsymbol{\varepsilon}_i$  is the random part.

Model 2 is “simple regression” in which the deterministic part  $\boldsymbol{\mu}$  is not constant but is a function of the nonrandom variable  $x$ . The assumption here is that this function is differentiable and can, in the range of the variation of the data, be approximated by a linear function [Tin51, pp. 19–20]. I.e., each element of  $\mathbf{y}$  is a constant  $\alpha$  plus a constant multiple of the corresponding element of the nonrandom vector  $\mathbf{x}$  plus a random error term:  $\mathbf{y}_t = \alpha + x_t \beta + \boldsymbol{\varepsilon}_t$ ,  $t = 1, \dots, n$ . This can be written as

$$(14.1.2) \quad \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \alpha + \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \beta + \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \vdots \\ \boldsymbol{\varepsilon}_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \vdots \\ \boldsymbol{\varepsilon}_n \end{bmatrix}$$

If the systematic part of  $\mathbf{y}$  depends on more than one variable, then one needs multiple regression, model 3. Mathematically, multiple regression has the same form (14.1.3), but this time  $\mathbf{X}$  is *arbitrary* (except for the restriction that all its columns are linearly independent). Model 3 has Models 1 and 2 as special cases.

Multiple regression is also used to “correct for” disturbing influences. Let us explain. A functional relationship, which makes the systematic part of  $\mathbf{y}$  depend on some other variable  $x$  will usually only hold if other relevant influences are kept constant. If those other influences vary, then they may affect the form of this functional relation. For instance, the marginal propensity to consume may be affected by the interest rate, or the unemployment rate. This is why some econometricians (Hendry) advocate that one should start with an “encompassing” model with many explanatory variables and then narrow the specification down by hypothesis testing. Milton Friedman, by contrast, is very suspicious about multiple regressions, and argues in [FS91, pp. 48/9] against the encompassing approach.

Friedman does not give a theoretical argument but argues by an example from Chemistry. Perhaps one can say that the variations in the other influences may have more serious implications than just modifying the form of the functional relationship: they may destroy this functional relation altogether, i.e., prevent any systematic or predictable behavior.

	observed	unobserved
random	$\mathbf{y}$	$\boldsymbol{\varepsilon}$
nonrandom	$\mathbf{X}$	$\boldsymbol{\beta}, \sigma^2$

### 14.2. Ordinary Least Squares

In the model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where  $\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2 \mathbf{I})$ , the OLS-estimate  $\hat{\boldsymbol{\beta}}$  is defined to be that value  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$  which minimizes

$$(14.2.1) \quad \text{SSE} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}.$$

Problem 156 shows that in model 1, this principle yields the arithmetic mean.

**PROBLEM 191.** *2 points Prove that, if one predicts a random variable  $y$  by a constant  $a$ , the constant which gives the best MSE is  $a = E[y]$ , and the best MSE one can get is  $\text{var}[y]$ .*

**ANSWER.**  $E[(y - a)^2] = E[y^2] - 2aE[y] + a^2$ . Differentiate with respect to  $a$  and set zero to get  $a = E[y]$ . One can also differentiate first and then take expected value:  $E[2(y - a)] = 0$ .  $\square$

We will solve this minimization problem using the first-order conditions in vector notation. As a preparation, you should read the beginning of Appendix C about matrix differentiation and the connection between matrix differentiation and the Jacobian matrix of a vector function. All you need at this point is the two equations (C.1.6) and (C.1.7). The chain rule (C.1.23) is enlightening but not strictly necessary for the present derivation.

The matrix differentiation rules (C.1.6) and (C.1.7) allow us to differentiate (14.2.1) to get

$$(14.2.2) \quad \partial \text{SSE} / \partial \boldsymbol{\beta}^\top = -2\mathbf{y}^\top \mathbf{X} + 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}.$$

Transpose it (because it is notationally simpler to have a relationship between column vectors), set it zero while at the same time replacing  $\boldsymbol{\beta}$  by  $\hat{\boldsymbol{\beta}}$ , and divide by 2, to get the “normal equation”

$$(14.2.3) \quad \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}}.$$

Due to our assumption that all columns of  $\mathbf{X}$  are linearly independent,  $\mathbf{X}^\top \mathbf{X}$  has an inverse and one can premultiply both sides of (14.2.3) by  $(\mathbf{X}^\top \mathbf{X})^{-1}$ :

$$(14.2.4) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

If the columns of  $\mathbf{X}$  are not linearly independent, then (14.2.3) has more than one solution, and the normal equation is also in this case a necessary and sufficient condition for  $\hat{\boldsymbol{\beta}}$  to minimize the SSE (proof in Problem 194).

**PROBLEM 192.** *4 points Using the matrix differentiation rules*

$$(14.2.5) \quad \partial \mathbf{w}^\top \mathbf{x} / \partial \mathbf{x}^\top = \mathbf{w}^\top$$

$$(14.2.6) \quad \partial \mathbf{x}^\top \mathbf{M} \mathbf{x} / \partial \mathbf{x}^\top = 2\mathbf{x}^\top \mathbf{M}$$

for symmetric  $\mathbf{M}$ , compute the least-squares estimate  $\hat{\boldsymbol{\beta}}$  which minimizes

$$(14.2.7) \quad \text{SSE} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

You are allowed to assume that  $\mathbf{X}^\top \mathbf{X}$  has an inverse.

**ANSWER.** First you have to multiply out

$$(14.2.8) \quad (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}.$$

The matrix differentiation rules (14.2.5) and (14.2.6) allow us to differentiate (14.2.8) to get

$$(14.2.9) \quad \partial \text{SSE} / \partial \boldsymbol{\beta}^\top = -2\mathbf{y}^\top \mathbf{X} + 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}.$$

Transpose it (because it is notationally simpler to have a relationship between column vectors), set it zero while at the same time replacing  $\boldsymbol{\beta}$  by  $\hat{\boldsymbol{\beta}}$ , and divide by 2, to get the “normal equation”

$$(14.2.10) \quad \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}}.$$

Since  $\mathbf{X}^\top \mathbf{X}$  has an inverse, one can premultiply both sides of (14.2.10) by  $(\mathbf{X}^\top \mathbf{X})^{-1}$ :

$$(14.2.11) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

**PROBLEM 193.** *2 points Show the following: if the columns of  $\mathbf{X}$  are linearly independent, then  $\mathbf{X}^\top \mathbf{X}$  has an inverse. ( $\mathbf{X}$  itself is not necessarily square.) In your proof you may use the following criteria: the columns of  $\mathbf{X}$  are linearly independent (this is also called:  $\mathbf{X}$  has full column rank) if and only if  $\mathbf{X}\mathbf{a} = \mathbf{o}$  implies  $\mathbf{a} = \mathbf{o}$ . And a square matrix has an inverse if and only if its columns are linearly independent.*

**ANSWER.** We have to show that any  $\mathbf{a}$  which satisfies  $\mathbf{X}^\top \mathbf{X}\mathbf{a} = \mathbf{o}$  is itself the null vector. From  $\mathbf{X}^\top \mathbf{X}\mathbf{a} = \mathbf{o}$  follows  $\mathbf{a}^\top \mathbf{X}^\top \mathbf{X}\mathbf{a} = 0$  which can also be written  $\|\mathbf{X}\mathbf{a}\|^2 = 0$ . Therefore  $\mathbf{X}\mathbf{a} = \mathbf{o}$  and since the columns of  $\mathbf{X}$  are linearly independent, this implies  $\mathbf{a} = \mathbf{o}$ .

**PROBLEM 194.** *3 points In this Problem we do not assume that  $\mathbf{X}$  has full column rank, it may be arbitrary.*

• a. *The normal equation (14.2.3) has always at least one solution. Hint: you are allowed to use, without proof, equation (A.3.3) in the mathematical appendix.*

**ANSWER.** With this hint it is easy:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  is a solution.

• b. *If  $\hat{\boldsymbol{\beta}}$  satisfies the normal equation and  $\boldsymbol{\beta}$  is an arbitrary vector, then*

$$(14.2.12) \quad (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})$$

**ANSWER.** This is true even if  $\mathbf{X}$  has deficient rank, and it will be shown here in this general case. To prove (14.2.12), write (14.2.1) as  $\text{SSE} = ((\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) - \mathbf{X}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}))^\top ((\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) - \mathbf{X}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}))$  since  $\hat{\boldsymbol{\beta}}$  satisfies (14.2.3), the cross product terms disappear.

• c. *Conclude from this that the normal equation is a necessary and sufficient condition characterizing the values  $\hat{\boldsymbol{\beta}}$  minimizing the sum of squared errors (14.2.1)*

ANSWER. (14.2.12) shows that the normal equations are sufficient. For necessity of the normal equations let  $\hat{\beta}$  be an arbitrary solution of the normal equation, we have seen that there is always at least one. Given  $\hat{\beta}$ , it follows from (14.2.12) that for any solution  $\beta^*$  of the minimization,  $\mathbf{X}^\top \mathbf{X}(\beta^* - \hat{\beta}) = \mathbf{o}$ . Use (14.2.3) to replace  $(\mathbf{X}^\top \mathbf{X})\hat{\beta}$  by  $\mathbf{X}^\top \mathbf{y}$  to get  $\mathbf{X}^\top \mathbf{X}\beta^* = \mathbf{X}^\top \mathbf{y}$ .  $\square$

It is customary to use the notation  $\mathbf{X}\hat{\beta} = \hat{\mathbf{y}}$  for the so-called *fitted values*, which are the estimates of the vector of means  $\boldsymbol{\eta} = \mathbf{X}\beta$ . Geometrically,  $\hat{\mathbf{y}}$  is the orthogonal projection of  $\mathbf{y}$  on the space spanned by the columns of  $\mathbf{X}$ . See Theorem A.6.1 about projection matrices.

The vector of differences between the actual and the fitted values is called the vector of “residuals”  $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \hat{\mathbf{y}}$ . The residuals are “predictors” of the actual (but unobserved) values of the disturbance vector  $\boldsymbol{\varepsilon}$ . An estimator of a random magnitude is usually called a “predictor,” but in the linear model estimation and prediction are treated on the same footing, therefore it is not necessary to distinguish between the two.

You should understand the difference between disturbances and residuals, and between the two decompositions

$$(14.2.13) \quad \mathbf{y} = \mathbf{X}\beta + \boldsymbol{\varepsilon} = \mathbf{X}\hat{\beta} + \hat{\boldsymbol{\varepsilon}}$$

PROBLEM 195. 2 points Assume that  $\mathbf{X}$  has full column rank. Show that  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}\mathbf{y}$  where  $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ . Show that  $\mathbf{M}$  is symmetric and idempotent.

ANSWER. By definition,  $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\beta} = \mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{y}$ . Idempotent, i.e.  $\mathbf{M}\mathbf{M} = \mathbf{M}$ :

$$(14.2.14) \quad \mathbf{M}\mathbf{M} = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)(\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{M} \quad \square$$

PROBLEM 196. Assume  $\mathbf{X}$  has full column rank. Define  $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .

• a. 1 point Show that the space  $\mathbf{M}$  projects on is the space orthogonal to all columns in  $\mathbf{X}$ , i.e.,  $\mathbf{M}\mathbf{q} = \mathbf{q}$  if and only if  $\mathbf{X}^\top \mathbf{q} = \mathbf{o}$ .

ANSWER.  $\mathbf{X}^\top \mathbf{q} = \mathbf{o}$  clearly implies  $\mathbf{M}\mathbf{q} = \mathbf{q}$ . Conversely,  $\mathbf{M}\mathbf{q} = \mathbf{q}$  implies  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{q} = \mathbf{o}$ . Premultiply this by  $\mathbf{X}^\top$  to get  $\mathbf{X}^\top \mathbf{q} = \mathbf{o}$ .  $\square$

• b. 1 point Show that a vector  $\mathbf{q}$  lies in the range space of  $\mathbf{X}$ , i.e., the space spanned by the columns of  $\mathbf{X}$ , if and only if  $\mathbf{M}\mathbf{q} = \mathbf{o}$ . In other words,  $\{\mathbf{q} : \mathbf{q} = \mathbf{X}\mathbf{a} \text{ for some } \mathbf{a}\} = \{\mathbf{q} : \mathbf{M}\mathbf{q} = \mathbf{o}\}$ .

ANSWER. First assume  $\mathbf{M}\mathbf{q} = \mathbf{o}$ . This means  $\mathbf{q} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{q} = \mathbf{X}\mathbf{a}$  with  $\mathbf{a} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{q}$ . Conversely, if  $\mathbf{q} = \mathbf{X}\mathbf{a}$  then  $\mathbf{M}\mathbf{q} = \mathbf{M}\mathbf{X}\mathbf{a} = \mathbf{O}\mathbf{a} = \mathbf{o}$ .  $\square$

PROBLEM 197. In 2-dimensional space, write down the projection matrix on diagonal line  $y = x$  (call it  $\mathbf{E}$ ), and compute  $\mathbf{E}\mathbf{z}$  for the three vectors  $\mathbf{a} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ,  $\mathbf{b} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ , and  $\mathbf{c} = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$ . Draw these vectors and their projections.

Assume we have a dependent variable  $\mathbf{y}$  and two regressors  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , each with 15 observations. Then one can visualize the data either as 15 points in 3-dimensional space (a 3-dimensional scatter plot), or 3 points in 15-dimensional space. In the first case, each point corresponds to an observation, in the second case, each point corresponds to a variable. In this latter case the points are usually represented as vectors. You only have 3 vectors, but each of these vectors is a vector in 15-dimensional space. But you do not have to draw a 15-dimensional space to draw these vectors; these 3 vectors span a 3-dimensional subspace, and  $\hat{\mathbf{y}}$  is the projection of the vector  $\mathbf{y}$  on the space spanned by the two regressors not only in the original 15-dimensional space, but already in this 3-dimensional subspace. In other words, [DM93, Figure 1.3] is valid in all dimensions! In the 15-dimensional space, each dimension represents one observation. In the 3-dimensional subspace, this is no longer true.

PROBLEM 198. “Simple regression” is regression with an intercept and one explanatory variable only, i.e.,

$$(14.2.15) \quad \mathbf{y}_t = \alpha + \beta x_t + \varepsilon_t$$

Here  $\mathbf{X} = \begin{bmatrix} 1 & \mathbf{x} \end{bmatrix}$  and  $\boldsymbol{\beta} = [\alpha \quad \beta]^\top$ . Evaluate (14.2.4) to get the following formulas for  $\hat{\boldsymbol{\beta}} = [\hat{\alpha} \quad \hat{\beta}]^\top$ :

$$(14.2.16) \quad \hat{\alpha} = \frac{\sum x_t^2 \sum \mathbf{y}_t - \sum x_t \sum x_t \mathbf{y}_t}{n \sum x_t^2 - (\sum x_t)^2}$$

$$(14.2.17) \quad \hat{\beta} = \frac{n \sum x_t \mathbf{y}_t - \sum x_t \sum \mathbf{y}_t}{n \sum x_t^2 - (\sum x_t)^2}$$

ANSWER.

$$(14.2.18) \quad \mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 1^\top \\ \mathbf{x}^\top \end{bmatrix} \begin{bmatrix} 1 & \mathbf{x} \end{bmatrix} = \begin{bmatrix} 1^\top 1 & 1^\top \mathbf{x} \\ \mathbf{x}^\top 1 & \mathbf{x}^\top \mathbf{x} \end{bmatrix} = \begin{bmatrix} n & \sum x_t \\ \sum x_t & \sum x_t^2 \end{bmatrix}$$

$$(14.2.19) \quad \mathbf{X}^\top \mathbf{X}^{-1} = \frac{1}{n \sum x_t^2 - (\sum x_t)^2} \begin{bmatrix} \sum x_t^2 & -\sum x_t \\ -\sum x_t & n \end{bmatrix}$$

$$(14.2.20) \quad \mathbf{X}^\top \mathbf{y} = \begin{bmatrix} 1^\top \mathbf{y} \\ \mathbf{x}^\top \mathbf{y} \end{bmatrix} = \begin{bmatrix} \sum \mathbf{y}_t \\ \sum x_i \mathbf{y}_t \end{bmatrix}$$

Therefore  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  gives equations (14.2.16) and (14.2.17).

PROBLEM 199. Show that

$$(14.2.21) \quad \sum_{t=1}^n (x_t - \bar{x})(y_t - \bar{y}) = \sum_{t=1}^n x_t y_t - n\bar{x}\bar{y}$$

(Note, as explained in [DM93, pp. 27/8] or [Gre97, Section 5.4.1], that the left hand side is computationally much more stable than the right.)

ANSWER. Simply multiply out.  $\square$

PROBLEM 200. Show that (14.2.17) and (14.2.16) can also be written as follows:

$$(14.2.22) \quad \hat{\beta} = \frac{\sum (x_t - \bar{x})(y_t - \bar{y})}{\sum (x_t - \bar{x})^2}$$

$$(14.2.23) \quad \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

ANSWER. Using  $\sum x_i = n\bar{x}$  and  $\sum y_i = n\bar{y}$  in (14.2.17), it can be written as

$$(14.2.24) \quad \hat{\beta} = \frac{\sum x_t y_t - n\bar{x}\bar{y}}{\sum x_t^2 - n\bar{x}^2}$$

Now apply Problem 199 to the numerator of (14.2.24), and Problem 199 with  $\mathbf{y} = \mathbf{x}$  to the denominator, to get (14.2.22).

To prove equation (14.2.23) for  $\hat{\alpha}$ , let us work backwards and plug (14.2.24) into the righthand side of (14.2.23):

$$(14.2.25) \quad \bar{y} - \bar{x}\hat{\beta} = \frac{\bar{y} \sum x_t^2 - \bar{y}n\bar{x}^2 - \bar{x} \sum x_t y_t + n\bar{x}\bar{x}\bar{y}}{\sum x_t^2 - n\bar{x}^2}$$

The second and the fourth term in the numerator cancel out, and what remains can be shown to be equal to (14.2.16).  $\square$

PROBLEM 201. 3 points Show that in the simple regression model, the fitted regression line can be written in the form

$$(14.2.26) \quad \hat{y}_t = \bar{y} + \hat{\beta}(x_t - \bar{x}).$$

From this follows in particular that the fitted regression line always goes through the point  $\bar{x}, \bar{y}$ .

ANSWER. Follows immediately if one plugs (14.2.23) into the defining equation  $\hat{y}_t = \hat{\alpha} + \hat{\beta}x_t$ .  $\square$

Formulas (14.2.22) and (14.2.23) are interesting because they express the regression coefficients in terms of the sample means and covariances. Problem 202 derives the properties of the population equivalents of these formulas:

PROBLEM 202. Given two random variables  $\mathbf{x}$  and  $\mathbf{y}$  with finite variances,  $\text{var}[\mathbf{x}] > 0$ . You know the expected values, variances and covariance of  $\mathbf{x}$  and  $\mathbf{y}$ ,  $\text{cov}[\mathbf{x}, \mathbf{y}]$ . You observe  $\mathbf{x}$ , but  $\mathbf{y}$  is unobserved. This question explores the properties of the Best Linear Unbiased Predictor (BLUP) of  $\mathbf{y}$  in this situation.

• a. 4 points Give a direct proof of the following, which is a special case of theorem 20.1.1: If you want to predict  $\mathbf{y}$  by an affine expression of the form  $a + b\mathbf{x}$ , you will get the lowest mean squared error MSE with  $b = \text{cov}[\mathbf{x}, \mathbf{y}] / \text{var}[\mathbf{x}]$  and  $a = \text{E}[\mathbf{y}] - b \text{E}[\mathbf{x}]$ .

ANSWER. The MSE is variance plus squared bias (see e.g. problem 165), therefore

$$(14.2.27) \quad \text{MSE}[a + b\mathbf{x}; \mathbf{y}] = \text{var}[a + b\mathbf{x} - \mathbf{y}] + (\text{E}[a + b\mathbf{x} - \mathbf{y}])^2 = \text{var}[b\mathbf{x} - \mathbf{y}] + (a - \text{E}[\mathbf{y}] + b \text{E}[\mathbf{x}])^2$$

Therefore we choose  $a$  so that the second term is zero, and then you only have to minimize the first term with respect to  $b$ . Since

$$(14.2.28) \quad \text{var}[b\mathbf{x} - \mathbf{y}] = b^2 \text{var}[\mathbf{x}] - 2b \text{cov}[\mathbf{x}, \mathbf{y}] + \text{var}[\mathbf{y}]$$

the first order condition is

$$(14.2.29) \quad 2b \text{var}[\mathbf{x}] - 2 \text{cov}[\mathbf{x}, \mathbf{y}] = 0$$

• b. 2 points For the first-order conditions you needed the partial derivatives  $\frac{\partial}{\partial a} \text{E}[(\mathbf{y} - a - b\mathbf{x})^2]$  and  $\frac{\partial}{\partial b} \text{E}[(\mathbf{y} - a - b\mathbf{x})^2]$ . It is also possible, and probably shorter, to interchange taking expected value and partial derivative, i.e., to compute  $\text{E}\left[\frac{\partial}{\partial a}(\mathbf{y} - a - b\mathbf{x})^2\right]$  and  $\text{E}\left[\frac{\partial}{\partial b}(\mathbf{y} - a - b\mathbf{x})^2\right]$  and set those zero. Do the above proof in this alternative fashion.

ANSWER.  $\text{E}\left[\frac{\partial}{\partial a}(\mathbf{y} - a - b\mathbf{x})^2\right] = -2 \text{E}[\mathbf{y} - a - b\mathbf{x}] = -2(\text{E}[\mathbf{y}] - a - b \text{E}[\mathbf{x}])$ . Setting this zero gives the formula for  $a$ . Now  $\text{E}\left[\frac{\partial}{\partial b}(\mathbf{y} - a - b\mathbf{x})^2\right] = -2 \text{E}[\mathbf{x}(\mathbf{y} - a - b\mathbf{x})] = -2(\text{E}[\mathbf{x}\mathbf{y}] - a \text{E}[\mathbf{x}] - b \text{E}[\mathbf{x}^2])$ . Setting this zero gives  $\text{E}[\mathbf{x}\mathbf{y}] - a \text{E}[\mathbf{x}] - b \text{E}[\mathbf{x}^2] = 0$ . Plug in formula for  $a$  and solve for  $b$ :

$$(14.2.30) \quad b = \frac{\text{E}[\mathbf{x}\mathbf{y}] - \text{E}[\mathbf{x}]\text{E}[\mathbf{y}]}{\text{E}[\mathbf{x}^2] - (\text{E}[\mathbf{x}])^2} = \frac{\text{cov}[\mathbf{x}, \mathbf{y}]}{\text{var}[\mathbf{x}]}.$$

• c. 2 points Compute the MSE of this predictor.

ANSWER. If one plugs the optimal  $a$  into (14.2.27), this just annuls the last term of (14.2.27) so that the MSE is given by (14.2.28). If one plugs the optimal  $b = \text{cov}[\mathbf{x}, \mathbf{y}] / \text{var}[\mathbf{x}]$  into (14.2.28) one gets

$$(14.2.31) \quad \text{MSE} = \left(\frac{\text{cov}[\mathbf{x}, \mathbf{y}]}{\text{var}[\mathbf{x}]}\right)^2 \text{var}[\mathbf{x}] - 2 \frac{\text{cov}[\mathbf{x}, \mathbf{y}]}{\text{var}[\mathbf{x}]} \text{cov}[\mathbf{x}, \mathbf{y}] + \text{var}[\mathbf{y}]$$

$$(14.2.32) \quad = \text{var}[\mathbf{y}] - \frac{(\text{cov}[\mathbf{x}, \mathbf{y}])^2}{\text{var}[\mathbf{x}]}.$$

- d. 2 points Show that the prediction error is uncorrelated with the observed  $x$ .

ANSWER.

$$(14.2.33) \quad \text{cov}[x, y - a - bx] = \text{cov}[x, y] - a \text{cov}[x, x] = 0$$

□

- e. 4 points If  $\text{var}[x] = 0$ , the quotient  $\text{cov}[x, y]/\text{var}[x]$  can no longer be formed, but if you replace the inverse by the  $g$ -inverse, so that the above formula becomes

$$(14.2.34) \quad b = \text{cov}[x, y](\text{var}[x])^{-}$$

then it always gives the minimum MSE predictor, whether or not  $\text{var}[x] = 0$ , and regardless of which  $g$ -inverse you use (in case there are more than one). To prove this, you need to answer the following four questions: (a) what is the BLUP if  $\text{var}[x] = 0$ ? (b) what is the  $g$ -inverse of a nonzero scalar? (c) what is the  $g$ -inverse of the scalar number 0? (d) if  $\text{var}[x] = 0$ , what do we know about  $\text{cov}[x, y]$ ?

ANSWER. (a) If  $\text{var}[x] = 0$  then  $x = \mu$  almost surely, therefore the observation of  $x$  does not give us any new information. The BLUP of  $y$  is  $\nu$  in this case, i.e., the above formula holds with  $b = 0$ .

(b) The  $g$ -inverse of a nonzero scalar is simply its inverse.

(c) Every scalar is a  $g$ -inverse of the scalar 0.

(d) if  $\text{var}[x] = 0$ , then  $\text{cov}[x, y] = 0$ .

Therefore pick a  $g$ -inverse 0, an arbitrary number will do, call it  $c$ . Then formula (14.2.34) says  $b = 0 \cdot c = 0$ . □

PROBLEM 203. 3 points Carefully state the specifications of the random variables involved in the linear regression model. How does the model in Problem 202 differ from the linear regression model? What do they have in common?

ANSWER. In the regression model, you have several observations, in the other model only one. In the regression model, the  $x_i$  are nonrandom, only the  $y_i$  are random, in the other model both  $x$  and  $y$  are random. In the regression model, the expected value of the  $y_i$  are not fully known, in the other model the expected values of both  $x$  and  $y$  are fully known. Both models have in common that the second moments are known only up to an unknown factor. Both models have in common that only first and second moments need to be known, and that they restrict themselves to linear estimators, and that the criterion function is the MSE (the regression model minimizes it, but the other model minimizes it since there is no unknown parameter whose value one has to minimize over. But this I cannot say right now, for this we need the Gauss-Markov theorem. Also the Gauss-Markov is valid in both cases!) □

PROBLEM 204. 2 points We are in the multiple regression model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with intercept, i.e.,  $\mathbf{X}$  is such that there is a vector  $\mathbf{a}$  with  $\boldsymbol{\iota} = \mathbf{X}\mathbf{a}$ . Define the row vector  $\bar{\mathbf{x}}^\top = \frac{1}{n}\boldsymbol{\iota}^\top \mathbf{X}$ , i.e., it has as its  $j$ th component the sample mean of the  $j$ th independent variable. Using the normal equations  $\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X}\hat{\boldsymbol{\beta}}$ , show that

$\bar{\mathbf{y}} = \bar{\mathbf{x}}^\top \hat{\boldsymbol{\beta}}$  (i.e., the regression plane goes through the center of gravity of all data points).

ANSWER. Premultiply the normal equation by  $\boldsymbol{\iota}^\top$  to get  $\boldsymbol{\iota}^\top \mathbf{y} - \boldsymbol{\iota}^\top \mathbf{X}\hat{\boldsymbol{\beta}} = 0$ . Premultiply by  $1/n$  to get the result.

PROBLEM 205. The fitted values  $\hat{\mathbf{y}}$  and the residuals  $\hat{\boldsymbol{\varepsilon}}$  are “orthogonal” in two different ways.

- a. 2 points Show that the inner product  $\hat{\mathbf{y}}^\top \hat{\boldsymbol{\varepsilon}} = 0$ . Why should you expect this from the geometric intuition of Least Squares?

ANSWER. Use  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}\mathbf{y}$  and  $\hat{\mathbf{y}} = (\mathbf{I} - \mathbf{M})\mathbf{y}$ :  $\hat{\mathbf{y}}^\top \hat{\boldsymbol{\varepsilon}} = \mathbf{y}^\top (\mathbf{I} - \mathbf{M})\mathbf{M}\mathbf{y} = 0$  because  $\mathbf{M}(\mathbf{I} - \mathbf{M}) = \mathbf{O}$ . This is a consequence of the more general result given in problem ??.

- b. 2 points Sometimes two random variables are called “orthogonal” if their covariance is zero. Show that  $\hat{\mathbf{y}}$  and  $\hat{\boldsymbol{\varepsilon}}$  are orthogonal also in this sense, i.e., show that for every  $i$  and  $j$ ,  $\text{cov}[\hat{y}_i, \hat{\varepsilon}_j] = 0$ . In matrix notation this can also be written  $\mathcal{C}[\hat{\mathbf{y}}, \hat{\boldsymbol{\varepsilon}}] = \mathbf{O}$ .

ANSWER.  $\mathcal{C}[\hat{\mathbf{y}}, \hat{\boldsymbol{\varepsilon}}] = \mathcal{C}[(\mathbf{I} - \mathbf{M})\mathbf{y}, \mathbf{M}\mathbf{y}] = (\mathbf{I} - \mathbf{M})\mathcal{V}[\mathbf{y}]\mathbf{M}^\top = (\mathbf{I} - \mathbf{M})(\sigma^2 \mathbf{I})\mathbf{M} = \sigma^2(\mathbf{I} - \mathbf{M})\mathbf{M} = \mathbf{O}$ . This is a consequence of the more general result given in question 246.

### 14.3. The Coefficient of Determination

Among the criteria which are often used to judge whether the model is appropriate, we will look at the “coefficient of determination”  $R^2$ , the “adjusted”  $\bar{R}^2$ , and later also at Mallows’  $C_p$  statistic. Mallows’  $C_p$  comes later because it is not a fit criterion but an initial criterion, i.e., it does not measure the fit of the model to the given data, but it estimates its MSE. Let us first look at  $R^2$ .

A value of  $R^2$  always is based (explicitly or implicitly) on a comparison of two models, usually nested in the sense that the model with fewer parameters can be viewed as a specialization of the model with more parameters. The value of  $R^2$  is then 1 minus the ratio of the smaller to the larger sum of squared residuals.

Thus, there is no such thing as the  $R^2$  from a single fitted model—one must always think about what model (perhaps an implicit “null” model) is held out as a standard of comparison. Once that is determined, the calculation is straightforward based on the sums of squared residuals from the two models. This is particularly appropriate for `nls()`, which minimizes a sum of squares.

The treatment which follows here is a little more complete than most. So far, textbooks, such as [DM93], never even give the leftmost term in formula (14.3.1) according to which  $R^2$  is the sample correlation coefficient. Other textbooks, such as [JHG+88] and [Gre97], do give this formula, but it remains a surprise: there is no explanation why the same quantity  $R^2$  can be expressed mathematically

two quite different ways, each of which has a different interpretation. The present treatment explains this.

If the regression has a constant term, then the OLS estimate  $\hat{\beta}$  has a third optimality property (in addition to minimizing the SSE and being the BLUE): no other linear combination of the explanatory variables has a higher squared sample correlation with  $\mathbf{y}$  than  $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$ .

In the proof of this optimality property we will use the symmetric and idempotent projection matrix  $\mathbf{D} = \mathbf{I} - \frac{1}{n}\boldsymbol{\iota}\boldsymbol{\iota}^\top$ . Applied to any vector  $\mathbf{z}$ ,  $\mathbf{D}$  gives  $\mathbf{D}\mathbf{z} = \mathbf{z} - \boldsymbol{\iota}\bar{z}$ , which is  $\mathbf{z}$  with the mean taken out. Taking out the mean is therefore a projection, on the space orthogonal to  $\boldsymbol{\iota}$ . See Problem 161.

PROBLEM 206. In the *reggeom* visualization, see Problem 293, in which  $\mathbf{x}_1$  is the vector of ones, which are the vectors  $\mathbf{D}\mathbf{x}_2$  and  $\mathbf{D}\mathbf{y}$ ?

ANSWER.  $\mathbf{D}\mathbf{x}_2$  is *og*, the dark blue line starting at the origin, and  $\mathbf{D}\mathbf{y}$  is *cy*, the red line starting on  $\mathbf{x}_1$  and going up to the peak.  $\square$

As an additional mathematical tool we will need the Cauchy-Schwartz inequality for the vector product:

$$(14.3.1) \quad (\mathbf{u}^\top \mathbf{v})^2 \leq (\mathbf{u}^\top \mathbf{u})(\mathbf{v}^\top \mathbf{v})$$

PROBLEM 207. If  $\mathbf{Q}$  is any nonnegative definite matrix, show that also

$$(14.3.2) \quad (\mathbf{u}^\top \mathbf{Q}\mathbf{v})^2 \leq (\mathbf{u}^\top \mathbf{Q}\mathbf{u})(\mathbf{v}^\top \mathbf{Q}\mathbf{v}).$$

ANSWER. This follows from the fact that any nnd matrix  $\mathbf{Q}$  can be written in the form  $\mathbf{Q} = \mathbf{R}^\top \mathbf{R}$ .  $\square$

In order to prove that  $\hat{\mathbf{y}}$  has the highest squared sample correlation, take any vector  $\mathbf{c}$  and look at  $\tilde{\mathbf{y}} = \mathbf{X}\mathbf{c}$ . We will show that the sample correlation of  $\mathbf{y}$  with  $\tilde{\mathbf{y}}$  cannot be higher than that of  $\mathbf{y}$  with  $\hat{\mathbf{y}}$ . For this let us first compute the sample covariance. By (9.3.17),  $n$  times the sample covariance between  $\tilde{\mathbf{y}}$  and  $\mathbf{y}$  is

$$(14.3.3) \quad n \text{ times sample covariance}(\tilde{\mathbf{y}}, \mathbf{y}) = \tilde{\mathbf{y}}^\top \mathbf{D}\mathbf{y} = \mathbf{c}^\top \mathbf{X}^\top \mathbf{D}(\hat{\mathbf{y}} + \hat{\boldsymbol{\epsilon}}).$$

By Problem 208,  $\mathbf{D}\hat{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\epsilon}}$ , hence  $\mathbf{X}^\top \mathbf{D}\hat{\boldsymbol{\epsilon}} = \mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = \mathbf{o}$  (this last equality is equivalent to the Normal Equation (14.2.3)), therefore (14.3.3) becomes  $\tilde{\mathbf{y}}^\top \mathbf{D}\mathbf{y} = \hat{\mathbf{y}}^\top \mathbf{D}\hat{\mathbf{y}}$ . Together with (14.3.2) this gives

$$(14.3.4) \quad (n \text{ times sample covariance}(\tilde{\mathbf{y}}, \mathbf{y}))^2 = (\tilde{\mathbf{y}}^\top \mathbf{D}\hat{\mathbf{y}})^2 \leq (\tilde{\mathbf{y}}^\top \mathbf{D}\tilde{\mathbf{y}})(\hat{\mathbf{y}}^\top \mathbf{D}\hat{\mathbf{y}})$$

In order to get from  $n^2$  times the squared sample covariance to the squared sample correlation coefficient we have to divide it by  $n^2$  times the sample variances

of  $\tilde{\mathbf{y}}$  and of  $\mathbf{y}$ :

$$(14.3.5) \quad (\text{sample correlation}(\tilde{\mathbf{y}}, \mathbf{y}))^2 = \frac{(\tilde{\mathbf{y}}^\top \mathbf{D}\mathbf{y})^2}{(\tilde{\mathbf{y}}^\top \mathbf{D}\tilde{\mathbf{y}})(\mathbf{y}^\top \mathbf{D}\mathbf{y})} \leq \frac{\hat{\mathbf{y}}^\top \mathbf{D}\hat{\mathbf{y}}}{\mathbf{y}^\top \mathbf{D}\mathbf{y}} = \frac{\sum(\hat{y}_j - \bar{y})^2}{\sum(y_j - \bar{y})^2} = \frac{\sum(\hat{y}_j - \bar{y})^2}{\sum(y_j - \bar{y})^2}$$

For the rightmost equal sign in (14.3.5) we need Problem 209.

If  $\tilde{\mathbf{y}} = \hat{\mathbf{y}}$ , inequality (14.3.4) becomes an equality, and therefore also (14.3.5) becomes an equality throughout. This completes the proof that  $\hat{\mathbf{y}}$  has the highest possible squared sample correlation with  $\mathbf{y}$ , and gives at the same time two different formulas for the same entity

$$(14.3.6) \quad R^2 = \frac{(\sum(\hat{y}_j - \bar{y})(y_j - \bar{y}))^2}{\sum(\hat{y}_j - \bar{y})^2 \sum(y_j - \bar{y})^2} = \frac{\sum(\hat{y}_j - \bar{y})^2}{\sum(y_j - \bar{y})^2}.$$

PROBLEM 208. 1 point Show that, if  $\mathbf{X}$  contains a constant term, then  $\mathbf{D}\hat{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\epsilon}}$ . You are allowed to use the fact that  $\mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = \mathbf{o}$ , which is equivalent to the normal equation (14.2.3).

ANSWER. Since  $\mathbf{X}$  has a constant term, a vector  $\mathbf{a}$  exists such that  $\mathbf{X}\mathbf{a} = \boldsymbol{\iota}$ , therefore  $\boldsymbol{\iota}^\top \mathbf{a} = \mathbf{a}^\top \mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = \mathbf{a}^\top \mathbf{o} = 0$ . From  $\boldsymbol{\iota}^\top \hat{\boldsymbol{\epsilon}} = 0$  follows  $\mathbf{D}\hat{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\epsilon}}$ .

PROBLEM 209. 1 point Show that, if  $\mathbf{X}$  has a constant term, then  $\bar{\hat{y}} = \bar{y}$

ANSWER. Follows from  $0 = \boldsymbol{\iota}^\top \hat{\boldsymbol{\epsilon}} = \boldsymbol{\iota}^\top \mathbf{y} - \boldsymbol{\iota}^\top \hat{\mathbf{y}}$ . In the visualization, this is equivalent with the fact that both *ocb* and *ocy* are right angles.

PROBLEM 210. Instead of (14.3.6) one often sees the formula

$$(14.3.7) \quad \frac{(\sum(\hat{y}_j - \bar{y})(y_j - \bar{y}))^2}{\sum(\hat{y}_j - \bar{y})^2 \sum(y_j - \bar{y})^2} = \frac{\sum(\hat{y}_j - \bar{y})^2}{\sum(y_j - \bar{y})^2}.$$

Prove that they are equivalent. Which equation is better?

The denominator in the righthand side expression of (14.3.6),  $\sum(y_j - \bar{y})^2$ , usually called “SST,” the total (corrected) sum of squares. The numerator  $\sum(\hat{y}_j - \bar{y})^2$  is usually called “SSR,” the sum of squares “explained” by the regression. In order to understand SSR better, we will show next the famous “Analysis of Variance” identity  $\text{SST} = \text{SSR} + \text{SSE}$ .

PROBLEM 211. In the *reggeom* visualization, again with  $\mathbf{x}_1$  representing vector of ones, show that  $\text{SST} = \text{SSR} + \text{SSE}$ , and show that  $R^2 = \cos^2 \alpha$  where  $\alpha$  is the angle between two lines in this visualization. Which lines?

ANSWER.  $\hat{\boldsymbol{\epsilon}}$  is the *by*, the green line going up to the peak, and SSE is the squared length of  $\hat{\boldsymbol{\epsilon}}$ . SST is the squared length of  $\mathbf{y} - \boldsymbol{\iota}\bar{y}$ . Since  $\boldsymbol{\iota}\bar{y}$  is the projection of  $\mathbf{y}$  on  $\mathbf{x}_1$ , i.e., it is *oc*, the part of  $\mathbf{x}_1$  that is red, one sees that SST is the squared length of *cy*. SSR is the squared length of *cb*.  $\square$



analysis of variance identity follows because  $cb\mathbf{y}$  is a right angle.  $R^2 = \cos^2 \alpha$  where  $\alpha$  is the angle between  $bc\mathbf{y}$  in this same triangle. □

Since the regression has a constant term, the decomposition

$$(14.3.8) \quad \mathbf{y} = (\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y}) + \boldsymbol{\nu}\bar{y}$$

is an orthogonal decomposition (all three vectors on the righthand side are orthogonal to each other), therefore in particular

$$(14.3.9) \quad (\mathbf{y} - \hat{\mathbf{y}})^\top (\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y}) = 0.$$

Geometrically this follows from the fact that  $\mathbf{y} - \hat{\mathbf{y}}$  is orthogonal to the column space of  $\mathbf{X}$ , while  $\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y}$  lies in that column space.

PROBLEM 212. *Show the decomposition 14.3.8 in the reggeom-visualization.*

ANSWER. From  $y$  take the green line down to  $b$ , then the light blue line to  $c$ , then the red line to the origin. □

This orthogonality can also be explained in terms of sequential projections: instead of projecting  $\mathbf{y}$  on  $\mathbf{x}_1$  directly I can first project it on the plane spanned by  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , and then project this projection on  $\mathbf{x}_1$ .

From (14.3.9) follows (now the same identity written in three different notations):

$$(14.3.10) \quad (\mathbf{y} - \boldsymbol{\nu}\bar{y})^\top (\mathbf{y} - \boldsymbol{\nu}\bar{y}) = (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y})^\top (\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y})$$

$$(14.3.11) \quad \sum_t (y_t - \bar{y})^2 = \sum_t (y_t - \hat{y}_t)^2 + \sum_t (\hat{y}_t - \bar{y})^2$$

$$(14.3.12) \quad \text{SST} = \text{SSE} + \text{SSR}$$

PROBLEM 213. *5 points Show that the “analysis of variance” identity  $\text{SST} = \text{SSE} + \text{SSR}$  holds in a regression with intercept, i.e., prove one of the two following equations:*

$$(14.3.13) \quad (\mathbf{y} - \boldsymbol{\nu}\bar{y})^\top (\mathbf{y} - \boldsymbol{\nu}\bar{y}) = (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y})^\top (\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y})$$

$$(14.3.14) \quad \sum_t (y_t - \bar{y})^2 = \sum_t (y_t - \hat{y}_t)^2 + \sum_t (\hat{y}_t - \bar{y})^2$$

ANSWER. Start with

$$(14.3.15) \quad \text{SST} = \sum (y_t - \bar{y})^2 = \sum (y_t - \hat{y}_t + \hat{y}_t - \bar{y})^2$$

and then show that the cross product term  $\sum (y_t - \hat{y}_t)(\hat{y}_t - \bar{y}) = \sum \hat{\varepsilon}_t(\hat{y}_t - \bar{y}) = \hat{\boldsymbol{\varepsilon}}^\top (\mathbf{X}\hat{\boldsymbol{\beta}} - \boldsymbol{\nu}\frac{1}{n}\boldsymbol{\nu}^\top \mathbf{y}) = 0$  since  $\hat{\boldsymbol{\varepsilon}}^\top \mathbf{X} = \mathbf{o}^\top$  and in particular, since a constant term is included,  $\hat{\boldsymbol{\varepsilon}}^\top \boldsymbol{\nu} = 0$ . □

From the so-called “analysis of variance” identity (14.3.12), together with (14.3.10) one obtains the following three alternative expressions for the maximum possible  $R^2$  relation, which is called  $R^2$  and which is routinely used as a measure of the “fit” of the regression:

$$(14.3.16) \quad R^2 = \frac{(\sum (\hat{y}_j - \bar{y})(y_j - \bar{y}))^2}{\sum (\hat{y}_j - \bar{y})^2 \sum (y_j - \bar{y})^2} = \frac{\text{SSR}}{\text{SST}} = \frac{\text{SST} - \text{SSE}}{\text{SST}}$$

The first of these three expressions is the squared sample correlation coefficient between  $\hat{\mathbf{y}}$  and  $\mathbf{y}$ , hence the notation  $R^2$ . The usual interpretation of the middle expression is the following: SST can be decomposed into a part SSR which is “explained” by the regression, and a part SSE which remains “unexplained,” and  $R^2$  measures that fraction of SST which can be “explained” by the regression. [Greiner pp. 250–253] and also [JHG+88, pp. 211/212] try to make this notion plausible. Instead of using the vague notions “explained” and “unexplained,” I prefer the following reading, which is based on the third expression for  $R^2$  in (14.3.16):  $\boldsymbol{\nu}\bar{y}$  is the vector of fitted values if one regresses  $\mathbf{y}$  on a constant term only, and SST is the SST in this “restricted” regression.  $R^2$  measures therefore the proportionate reduction in the SSE if one adds the nonconstant regressors to the regression. From this last formula one can also see that  $R^2 = \cos^2 \alpha$  where  $\alpha$  is the angle between  $\mathbf{y} - \boldsymbol{\nu}\bar{y}$  and  $\hat{\mathbf{y}} - \boldsymbol{\nu}\bar{y}$ .

PROBLEM 214. *Given two data series  $\mathbf{x}$  and  $\mathbf{y}$ . Show that the regression of  $\mathbf{y}$  on  $\mathbf{x}$  has the same  $R^2$  as the regression of  $\mathbf{x}$  on  $\mathbf{y}$ . (Both regressions are assumed to include a constant term.) Easy, but you have to think!*

ANSWER. The symmetry comes from the fact that, in this particular case,  $R^2$  is the squared sample correlation coefficient between  $\mathbf{x}$  and  $\mathbf{y}$ . Proof:  $\hat{\mathbf{y}}$  is an affine transformation of  $\mathbf{x}$ , and correlation coefficients are invariant under affine transformations (compare Problem 216).

PROBLEM 215. *This Problem derives some relationships which are valid in simple regression  $y_t = \alpha + \beta x_t + \varepsilon_t$  but their generalization to multiple regression is not obvious.*

• a. *2 points Show that*

$$(14.3.17) \quad R^2 = \hat{\beta}^2 \frac{\sum (x_t - \bar{x})^2}{\sum (y_t - \bar{y})^2}$$

*Hint: show first that  $\hat{y}_t - \bar{y} = \hat{\beta}(x_t - \bar{x})$ .*

ANSWER. From  $\hat{y}_t = \hat{\alpha} + \hat{\beta}x_t$  and  $\bar{y} = \hat{\alpha} + \hat{\beta}\bar{x}$  follows  $\hat{y}_t - \bar{y} = \hat{\beta}(x_t - \bar{x})$ . Therefore

$$(14.3.18) \quad R^2 = \frac{\sum (\hat{y}_t - \bar{y})^2}{\sum (y_t - \bar{y})^2} = \hat{\beta}^2 \frac{\sum (x_t - \bar{x})^2}{\sum (y_t - \bar{y})^2}$$

• b. 2 points Furthermore show that  $R^2$  is the sample correlation coefficient between  $\mathbf{y}$  and  $\mathbf{x}$ , i.e.,

$$(14.3.19) \quad R^2 = \frac{\left(\sum(x_t - \bar{x})(y_t - \bar{y})\right)^2}{\sum(x_t - \bar{x})^2 \sum(y_t - \bar{y})^2}.$$

Hint: you are allowed to use (14.2.22).

ANSWER.

$$(14.3.20) \quad R^2 = \hat{\beta}^2 \frac{\sum(x_t - \bar{x})^2}{\sum(y_t - \bar{y})^2} = \frac{\left(\sum(x_t - \bar{x})(y_t - \bar{y})\right)^2 \sum(x_t - \bar{x})^2}{\left(\sum(x_t - \bar{x})^2\right)^2 \sum(y_t - \bar{y})^2}$$

which simplifies to (14.3.19).  $\square$

• c. 1 point Finally show that  $R^2 = \hat{\beta}_{xy}\hat{\beta}_{yx}$ , i.e., it is the product of the two slope coefficients one gets if one regresses  $\mathbf{y}$  on  $\mathbf{x}$  and  $\mathbf{x}$  on  $\mathbf{y}$ .

If the regression does not have a constant term, but a vector  $\mathbf{a}$  exists with  $\boldsymbol{\iota} = \mathbf{X}\mathbf{a}$ , then the above mathematics remains valid. If  $\mathbf{a}$  does not exist, then the identity  $\text{SST} = \text{SSR} + \text{SSE}$  no longer holds, and (14.3.16) is no longer valid. The fraction  $\frac{\text{SST} - \text{SSE}}{\text{SST}}$  can assume negative values. Also the sample correlation coefficient between  $\hat{\mathbf{y}}$  and  $\mathbf{y}$  loses its motivation, since there will usually be other linear combinations of the columns of  $\mathbf{X}$  that have higher sample correlation with  $\mathbf{y}$  than the fitted values  $\hat{\mathbf{y}}$ .

Equation (14.3.16) is still puzzling at this point: why do two quite different simple concepts, the sample correlation and the proportionate reduction of the SSE, give the same numerical result? To explain this, we will take a short digression about correlation coefficients, in which it will be shown that correlation coefficients *always* denote proportionate reductions in the MSE. Since the SSE is (up to a constant factor) the sample equivalent of the MSE of the prediction of  $\mathbf{y}$  by  $\hat{\mathbf{y}}$ , this shows that (14.3.16) is simply the sample equivalent of a general fact about correlation coefficients.

But first let us take a brief look at the Adjusted  $R^2$ .

#### 14.4. The Adjusted R-Square

The coefficient of determination  $R^2$  is often used as a criterion for the selection of regressors. There are several drawbacks to this. [KA69, Chapter 8] shows that the distribution function of  $R^2$  depends on both the unknown error variance and the values taken by the explanatory variables; therefore the  $R^2$  belonging to different regressions cannot be compared.

A further drawback is that inclusion of more regressors always increases  $R^2$ . The adjusted  $\bar{R}^2$  is designed to remedy this. Starting from the formula  $R^2 = 1 - \text{SSE}/\text{SST}$ , the “adjustment” consists in dividing both SSE and SST by the degrees of freedom:

$$(14.4.1) \quad \bar{R}^2 = 1 - \frac{\text{SSE}/(n-k)}{\text{SST}/(n-1)} = 1 - (1 - R^2) \frac{n-1}{n-k}.$$

For given SST, i.e., when one looks at alternative regressions with the same dependent variable,  $\bar{R}^2$  is therefore a declining function of  $s^2$ , the unbiased estimator of  $\sigma^2$ . Choosing the regression with the highest  $\bar{R}^2$  amounts therefore to selecting the regression which yields the lowest value for  $s^2$ .

$\bar{R}^2$  has the following interesting property: (which we note here only for reference because we have not yet discussed the F-test:) Assume one adds  $i$  more regressors then  $\bar{R}^2$  increases only if the F statistic for these additional regressors has a value greater than one. One can also say:  $s^2$  decreases only if  $F > 1$ . To see this, write this F statistic as

$$(14.4.2) \quad F = \frac{(\text{SSE}_k - \text{SSE}_{k+i})/i}{\text{SSE}_{k+i}/(n-k-i)} = \frac{n-k-i}{i} \left( \frac{\text{SSE}_k}{\text{SSE}_{k+i}} - 1 \right)$$

$$(14.4.3) \quad = \frac{n-k-i}{i} \left( \frac{(n-k)s_k^2}{(n-k-i)s_{k+i}^2} - 1 \right)$$

$$(14.4.4) \quad = \frac{(n-k)s_k^2}{is_{k+i}^2} - \frac{n-k}{i} + 1$$

$$(14.4.5) \quad = \frac{(n-k)}{i} \left( \frac{s_k^2}{s_{k+i}^2} - 1 \right) + 1$$

From this the statement follows.

Minimizing the adjusted  $\bar{R}^2$  is equivalent to minimizing the unbiased variance estimator  $s^2$ ; it still does not penalize the loss of degrees of freedom heavily enough, i.e., it still admits too many variables into the model.

Alternatives minimize Amemiya’s prediction criterion or Akaike’s information criterion, which minimize functions of the estimated variances and  $n$  and  $k$ . Akaike’s information criterion minimizes an estimate of the Kullback-Leibler discrepancy which was discussed on p. 158.

## CHAPTER 15

## Digression about Correlation Coefficients

## 15.1. A Unified Definition of Correlation Coefficients

Correlation coefficients measure linear association. The usual definition of the simple correlation coefficient between two variables  $\rho_{xy}$  (sometimes we also use the notation  $\text{corr}[x, y]$ ) is their standardized covariance

$$(15.1.1) \quad \rho_{xy} = \frac{\text{cov}[x, y]}{\sqrt{\text{var}[x]}\sqrt{\text{var}[y]}}.$$

Because of Cauchy-Schwartz, its value lies between  $-1$  and  $1$ .

**PROBLEM 216.** *Given the constant scalars  $a \neq 0$  and  $c \neq 0$  and  $b$  and  $d$  arbitrary. Show that  $\text{corr}[x, y] = \pm \text{corr}[ax + b, cy + d]$ , with the  $+$  sign being valid if  $a$  and  $c$  have the same sign, and the  $-$  sign otherwise.*

**ANSWER.** Start with  $\text{cov}[ax + b, cy + d] = ac \text{cov}[x, y]$  and go from there.  $\square$

Besides the *simple* correlation coefficient  $\rho_{xy}$  between two scalar variables  $y$  and  $x$ , one can also define the squared *multiple* correlation coefficient  $\rho_{y(x)}^2$  between one scalar variable  $y$  and a whole vector of variables  $\mathbf{x}$ , and the *partial* correlation coefficient  $\rho_{12 \cdot \mathbf{x}}$  between two scalar variables  $y_1$  and  $y_2$ , with a vector of other variables  $\mathbf{x}$  “partialled out.” The multiple correlation coefficient measures the strength of a linear association between  $y$  and *all* components of  $\mathbf{x}$  together, and the partial correlation coefficient measures the strength of that part of the linear association between  $y_1$  and  $y_2$  which cannot be attributed to their joint association with  $\mathbf{x}$ . One can also define partial multiple correlation coefficients. If one wants to measure the linear association between two *vectors*, then one number is no longer enough, but one needs several numbers, the “canonical correlations.”

The multiple or partial correlation coefficients are usually defined as simple correlation coefficients involving the best linear predictor or its residual. But all these correlation coefficients share the property that they indicate a proportionate reduction in the MSE. See e.g. [Rao73, pp. 268–70]. Problem 217 makes this point for the *simple* correlation coefficient:

**PROBLEM 217.** *4 points Show that the proportionate reduction in the MSE of the best predictor of  $y$ , if one goes from predictors of the form  $y^* = a$  to predictors of the form  $y^* = a + bx$ , is equal to the squared correlation coefficient between  $y$  and  $x$ . You are allowed to use the results of Problems 191 and 202. To set notation, let  $\nu$  be the minimum MSE in the first prediction (Problem 191)  $\text{MSE}[\text{constant term}; y]$ , and  $\omega$  be the minimum MSE in the second prediction (Problem 202)  $\text{MSE}[\text{constant term and } x; y]$ . Show that*

$$(15.1.2) \quad \frac{\text{MSE}[\text{constant term}; y] - \text{MSE}[\text{constant term and } x; y]}{\text{MSE}[\text{constant term}; y]} = \frac{(\text{cov}[y, x])^2}{\text{var}[y] \text{var}[x]} = \rho_{yx}^2.$$

**ANSWER.** The minimum MSE with only a constant is  $\text{var}[y]$  and (14.2.32) says that  $\text{MSE}[\text{constant term and } x; y] = \text{var}[y] - (\text{cov}[x, y])^2 / \text{var}[x]$ . Therefore the difference in MSE’s is  $(\text{cov}[x, y])^2 / \text{var}[x]$  and if one divides by  $\text{var}[y]$  to get the relative difference, one gets exactly the squared correlation coefficient.

*Multiple Correlation Coefficients.* Now assume  $\mathbf{x}$  is a vector while  $y$  remains a scalar. Their joint mean vector and dispersion matrix are

$$(15.1.3) \quad \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{\mu} \\ \nu \end{bmatrix}, \sigma^2 \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}} & \boldsymbol{\omega}_{\mathbf{x}y} \\ \boldsymbol{\omega}_{\mathbf{x}y}^\top & \omega_{yy} \end{bmatrix}.$$

By theorem ??, the best linear predictor of  $y$  based on  $\mathbf{x}$  has the formula

$$(15.1.4) \quad y^* = \nu + \boldsymbol{\omega}_{\mathbf{x}y}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$y^*$  has the following additional extremal value property: no linear combination  $\mathbf{b}^\top \mathbf{x}$  has a higher squared correlation with  $y$  than  $y^*$ . This maximal value of the squared correlation is called the squared multiple correlation coefficient

$$(15.1.5) \quad \rho_{y(x)}^2 = \frac{\boldsymbol{\omega}_{\mathbf{x}y}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\omega}_{\mathbf{x}y}}{\omega_{yy}}.$$

The multiple correlation coefficient itself is the positive square root, i.e., it is always nonnegative, while some other correlation coefficients may take on negative values.

The squared multiple correlation coefficient can also be defined in terms of proportionate reduction in MSE. It is equal to the proportionate reduction in the MSE of the best predictor of  $y$  if one goes from predictors of the form  $y^* = a$  to predictors of the form  $y^* = a + \mathbf{b}^\top \mathbf{x}$ , i.e.,

$$(15.1.6) \quad \rho_{y(x)}^2 = \frac{\text{MSE}[\text{constant term}; y] - \text{MSE}[\text{constant term and } \mathbf{x}; y]}{\text{MSE}[\text{constant term}; y]}$$

There are therefore two natural definitions of the multiple correlation coefficient. These two definitions correspond to the two formulas for  $R^2$  in (14.3.6).

*Partial Correlation Coefficients.* Now assume  $\mathbf{y} = [\mathbf{y}_1 \ \mathbf{y}_2]^\top$  is a vector with two elements and write

$$(15.1.7) \quad \begin{bmatrix} \mathbf{x} \\ \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{\mu} \\ \nu_1 \\ \nu_2 \end{bmatrix}, \sigma^2 \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}} & \boldsymbol{\omega}_{\mathbf{y}_1} & \boldsymbol{\omega}_{\mathbf{y}_2} \\ \boldsymbol{\omega}_{\mathbf{y}_1}^\top & \omega_{11} & \omega_{12} \\ \boldsymbol{\omega}_{\mathbf{y}_2}^\top & \omega_{21} & \omega_{22} \end{bmatrix}.$$

Let  $\mathbf{y}^*$  be the best linear predictor of  $\mathbf{y}$  based on  $\mathbf{x}$ . The partial correlation coefficient  $\rho_{12.\mathbf{x}}$  is defined to be the simple correlation between the residuals  $\text{corr}[(\mathbf{y}_1 - \mathbf{y}_1^*), (\mathbf{y}_2 - \mathbf{y}_2^*)]$ . This measures the correlation between  $\mathbf{y}_1$  and  $\mathbf{y}_2$  which is “local,” i.e., which does not follow from their association with  $\mathbf{x}$ . Assume for instance that both  $\mathbf{y}_1$  and  $\mathbf{y}_2$  are highly correlated with  $\mathbf{x}$ . Then they will also have a high correlation with each other. Subtracting  $\mathbf{y}_i^*$  from  $\mathbf{y}_i$  eliminates this dependency on  $\mathbf{x}$ , therefore any remaining correlation is “local.” Compare [Krz88, p. 475].

The partial correlation coefficient can be defined as the relative reduction in the MSE if one adds  $\mathbf{y}_2$  to  $\mathbf{x}$  as a predictor of  $\mathbf{y}_1$ :

$$(15.1.8) \quad \rho_{12.\mathbf{x}}^2 = \frac{\text{MSE}[\text{constant term and } \mathbf{x}; \mathbf{y}_2] - \text{MSE}[\text{constant term, } \mathbf{x}, \text{ and } \mathbf{y}_1; \mathbf{y}_2]}{\text{MSE}[\text{constant term and } \mathbf{x}; \mathbf{y}_2]}.$$

PROBLEM 218. Using the definitions in terms of MSE’s, show that the following relationship holds between the squares of multiple and partial correlation coefficients:

$$(15.1.9) \quad 1 - \rho_{2(\mathbf{x},1)}^2 = (1 - \rho_{21.\mathbf{x}}^2)(1 - \rho_{2(\mathbf{x})}^2)$$

ANSWER. In terms of the MSE, (15.1.9) reads

$$(15.1.10) \quad \frac{\text{MSE}[\text{constant term, } \mathbf{x}, \text{ and } \mathbf{y}_1; \mathbf{y}_2]}{\text{MSE}[\text{constant term}; \mathbf{y}_2]} = \frac{\text{MSE}[\text{constant term, } \mathbf{x}, \text{ and } \mathbf{y}_1; \mathbf{y}_2]}{\text{MSE}[\text{constant term and } \mathbf{x}; \mathbf{y}_2]} \frac{\text{MSE}[\text{constant term and } \mathbf{x}; \mathbf{y}_2]}{\text{MSE}[\text{constant term}; \mathbf{y}_2]}$$

□

From (15.1.9) follows the following weighted average formula:

$$(15.1.11) \quad \rho_{2(\mathbf{x},1)}^2 = \rho_{2(\mathbf{x})}^2 + (1 - \rho_{2(\mathbf{x})}^2)\rho_{21.\mathbf{x}}^2$$

An alternative proof of (15.1.11) is given in [Gra76, pp. 116/17].

*Mixed cases:* One can also form multiple correlations coefficients with some of the variables partialled out. The dot notation used here is due to Yule, [Yul07]. The notation, definition, and formula for the squared correlation coefficient is

$$(15.1.12) \quad \rho_{\mathbf{y}(\mathbf{x}).\mathbf{z}}^2 = \frac{\text{MSE}[\text{constant term and } \mathbf{z}; \mathbf{y}] - \text{MSE}[\text{constant term, } \mathbf{z}, \text{ and } \mathbf{x}; \mathbf{y}]}{\text{MSE}[\text{constant term and } \mathbf{z}; \mathbf{y}]}$$

$$(15.1.13) \quad = \frac{\boldsymbol{\omega}_{\mathbf{x}\mathbf{y}.z}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}.z}^- \boldsymbol{\omega}_{\mathbf{x}\mathbf{y}.z}}{\omega_{\mathbf{y}\mathbf{y}.z}}$$

## CHAPTER 16

## Specific Datasets

## 16.1. Cobb Douglas Aggregate Production Function

PROBLEM 219. *2 points* The Cobb-Douglas production function postulates the following relationship between annual output  $q_t$  and the inputs of labor  $\ell_t$  and capital  $k_t$ :

$$(16.1.1) \quad q_t = \mu \ell_t^\beta k_t^\gamma \exp(\varepsilon_t).$$

$q_t$ ,  $\ell_t$ , and  $k_t$  are observed, and  $\mu$ ,  $\beta$ ,  $\gamma$ , and the  $\varepsilon_t$  are to be estimated. By the variable transformation  $x_t = \log q_t$ ,  $y_t = \log \ell_t$ ,  $z_t = \log k_t$ , and  $\alpha = \log \mu$ , one obtains the linear regression

$$(16.1.2) \quad x_t = \alpha + \beta y_t + \gamma z_t + \varepsilon_t$$

Sometimes the following alternative variable transformation is made:  $u_t = \log(q_t/\ell_t)$ ,  $v_t = \log(k_t/\ell_t)$ , and the regression

$$(16.1.3) \quad u_t = \alpha + \gamma v_t + \varepsilon_t$$

is estimated. How are the regressions (16.1.2) and (16.1.3) related to each other?

ANSWER. Write (16.1.3) as

$$(16.1.4) \quad x_t - y_t = \alpha + \gamma(z_t - y_t) + \varepsilon_t$$

and collect terms to get

$$(16.1.5) \quad x_t = \alpha + (1 - \gamma)y_t + \gamma z_t + \varepsilon_t$$

From this follows that running the regression (16.1.3) is equivalent to running the regression (16.1.2) with the constraint  $\beta + \gamma = 1$  imposed.  $\square$

The assumption here is that output is the only random variable. The regression model is based on the assumption that the dependent variables have more noise in them than the independent variables. One can justify this by the argument that any noise in the independent variables will be transferred to the dependent variable, and also that variables which affect other variables have more steadiness in them than variables which depend on others. This justification often has merit, but in the specific case, there is much more measurement error in the labor and capital inputs

than in the outputs. Therefore the assumption that only the output has an error term is clearly wrong, and problem 221 below will look for possible alternatives.

PROBLEM 220. *Table 1 shows the data used by Cobb and Douglas in their original article [CD28] introducing the production function which would bear their name. The output is “Day’s index of the physical volume of production (1899 = 100)” described in [DP20], capital is the capital stock in manufacturing in millions of 1880 dollars [CD28, p. 145], labor is the “probable average number of wage earners employed in manufacturing” [CD28, p. 148], and wage is an index of the real wage (1899–1919 = 100).*

year	1899	1900	1901	1902	1903	1904	1905	1906	1907	1908	1909	1910
output	100	101	112	122	124	122	143	152	151	126	155	159
capital	4449	4746	5061	5444	5806	6132	6626	7234	7832	8229	8820	9240
labor	4713	4968	5184	5554	5784	5468	5906	6251	6483	5714	6615	6807
wage	99	98	101	102	100	99	103	101	99	94	102	104
year	1911	1912	1913	1914	1915	1916	1917	1918	1919	1920	1921	1922
output	153	177	184	169	189	225	227	223	218	231	179	240
capital	9624	10067	10520	10873	11840	13242	14915	16265	17234	18118	18542	19192
labor	6855	7167	7277	7026	7269	8601	9218	9446	9096	9110	6947	7602
wage	97	99	100	99	99	104	103	107	111	114	115	119

TABLE 1. Cobb Douglas Original Data

• a. A text file with the data is available on the web at [www.econ.utah.edu/ehrbardata/cobbdoug.txt](http://www.econ.utah.edu/ehrbardata/cobbdoug.txt), and a SDML file (XML for statistical data which can be read by R, Matlab, and perhaps also SPSS) is available at [www.econ.utah.edu/ehrbardata/cobbdoug.sdml](http://www.econ.utah.edu/ehrbardata/cobbdoug.sdml). Load these data into your favorite statistics package.

ANSWER. In R, you can simply issue the command `cobbdoug <- read.table("http://www.econ.utah.edu/ehrbardata/cobbdoug.txt", header=TRUE)`. If you run R on unix, you can also do the following: download `cobbdoug.sdml` from the web, and then first issue the command `library(StatDataML)` and then `readSDML("cobbdoug.sdml")`. When I tried this last, the XML package necessary for `StatDataML` was not available on windows, but chances are it will be when you read this.

In SAS, you must issue the commands

```
data cobbdoug;
  infile 'cobbdoug.txt';
  input year output capital labor;
run;
```

But for this to work you must delete the first line in the file `cobbdoug.txt` which contains the variable names. (Is it possible to tell SAS to skip the first line?) And you may have to tell SAS

the full pathname of the text file with the data. If you want a permanent instead of a temporary dataset, give it a two-part name, such as `ecmet.cobbdoug`.

Here are the instructions for SPSS: 1) Begin SPSS with a blank spreadsheet. 2) Open up a file with the following commands and run:

```
SET
BLANKS=SYSMIS
UNDEFINED=WARN.
DATA LIST
FILE='A:\Cbbunst.dat' FIXED RECORDS=1 TABLE /1 year 1-4 output 5-9 capital
10-16 labor 17-22 wage 23-27 .
EXECUTE.
```

This files assume the data file to be on the same directory, and again the first line in the data file with the variable names must be deleted. Once the data are entered into SPSS the procedures (regression, etc.) are best run from the point and click environment.  $\square$

- b. *The next step is to look at the data. On [CD28, p. 150], Cobb and Douglas plot `capital`, `labor`, and `output` on a logarithmic scale against time, all 3 series normalized such that they start in 1899 at the same level =100. Reproduce this graph using a modern statistics package.*

- c. *Run both regressions (16.1.2) and (16.1.3) on Cobb and Douglas's original dataset. Compute 95% confidence intervals for the coefficients of capital and labor in the unconstrained and the constrained models.*

ANSWER. SAS does not allow you to transform the data on the fly, it insists that you first go through a data step creating the transformed data, before you can run a regression on them. Therefore the next set of commands creates a temporary dataset `cdtmp`. The data step `data cdtmp` includes all the data from `cobbdoug` into `cdtmp` and then creates some transformed data as well. Then one can run the regressions. Here are the commands; they are in the file `cbbrgrss.sas` in your data disk:

```
data cdtmp;
  set cobbdoug;
  logcap = log(capital);
  loglab = log(labor);
  logout = log(output);
  logcl = logcap-loglab;
  logol = logout-loglab;
run;
proc reg data = cdtmp;
  model logout = logcap loglab;
run;
proc reg data = cdtmp;
  model logol = logcl;
run;
```

Careful! In R, the command `lm(log(output)-log(labor) ~ log(capital)-log(labor), da` does not give the right results. It does not complain but the result is wrong nevertheless. The right way to write this command is `lm(I(log(output)-log(labor)) ~ I(log(capital)-log(labor)),`

- d. *The regression results are graphically represented in Figure 1. The ellipse is a joint 95% confidence region for  $\beta$  and  $\gamma$ . This ellipse is a level line of SSE. The vertical and horizontal bands represent univariate 95% confidence regions for  $\beta$  and  $\gamma$  separately. The diagonal line is the set of all  $\beta$  and  $\gamma$  with  $\beta + \gamma = 1$  representing the constraint of constant returns to scale. The small ellipse is that level line of the SSE which is tangent to the constraint. The point of tangency represents the constrained estimator. Reproduce this graph (or as much of this graph as you can) using your statistics package.*

Remark: In order to make the hand computations easier, Cobb and Douglas reduced the data for `capital` and `labor` to index numbers (1899=100) which were rounded to integers, before running the regressions, and Figure 1 was constructed using these rounded data. Since you are using the nonstandardized data, you may get slightly different results.

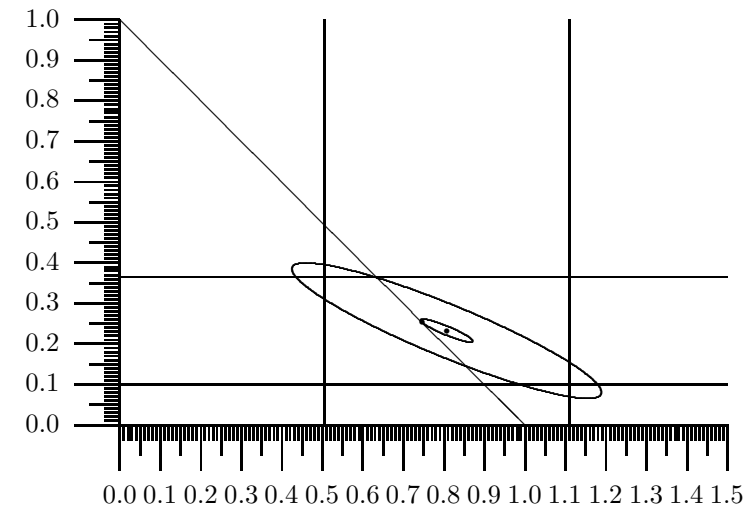


FIGURE 1. Coefficients of capital (vertical) and labor (horizontal), dependent variable output, unconstrained and constrained, 1899–1922

ANSWER. `lines(ellipse.lm(cbbfit, which=c(2, 3)))`

PROBLEM 221. *In this problem we will treat the Cobb-Douglas data as a dataset with errors in all three variables. See chapter ?? and problem ?? about that.*

- a. *Run the three elementary regressions for the whole period, then choose at least two subperiods and run it for those. Plot all regression coefficients as points in a plane, using different colors for the different subperiods (you have to normalize them in a special way that they all fit on the same plot).*

ANSWER. Here are the results in R:

```
> outputlm<-lm(log(output)~log(capital)+log(labor),data=cobbdoug)
> capitallm<-lm(log(capital)~log(labor)+log(output),data=cobbdoug)
> laborlm<-lm(log(labor)~log(output)+log(capital),data=cobbdoug)
> coefficients(outputlm)
(Intercept) log(capital) log(labor)
-0.1773097 0.2330535 0.8072782
> coefficients(capitallm)
(Intercept) log(labor) log(output)
-2.72052726 -0.08695944 1.67579357
> coefficients(laborlm)
(Intercept) log(output) log(capital)
1.27424214 0.73812541 -0.01105754

#Here is the information for the confidence ellipse:
> summary(outputlm,correlation=T)

Call:
lm(formula = log(output) ~ log(capital) + log(labor), data = cobbdoug)

Residuals:
    Min       1Q   Median       3Q      Max
-0.075282 -0.035234 -0.006439  0.038782  0.142114

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.17731  0.43429   -0.408  0.68721
log(capital)  0.23305  0.06353    3.668  0.00143 **
log(labor)   0.80728  0.14508    5.565  1.6e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.05814 on 21 degrees of freedom
Multiple R-Squared: 0.9574, Adjusted R-squared: 0.9534
F-statistic: 236.1 on 2 and 21 degrees of freedom, p-value: 3.997e-15

Correlation of Coefficients:
(Intercept) log(capital)
log(capital) 0.7243
```

```
log(labor) -0.9451 -0.9096

#Quantile of the F-distribution:
> qf(p=0.95, df1=2, df2=21)

[1] 3.4668
```

- b. *The elementary regressions will give you three fitted equations of the form*

$$(16.1.6) \quad output = \hat{\alpha}_1 + \hat{\beta}_{12} labor + \hat{\beta}_{13} capital + residual_1$$

$$(16.1.7) \quad labor = \hat{\alpha}_2 + \hat{\beta}_{21} output + \hat{\beta}_{23} capital + residual_2$$

$$(16.1.8) \quad capital = \hat{\alpha}_3 + \hat{\beta}_{31} output + \hat{\beta}_{32} labor + residual_3.$$

*In order to compare the slope parameters in these regressions, first rearrange them in the form*

$$(16.1.9) \quad -output + \hat{\beta}_{12} labor + \hat{\beta}_{13} capital + \hat{\alpha}_1 + residual_1 = 0$$

$$(16.1.10) \quad \hat{\beta}_{21} output - labor + \hat{\beta}_{23} capital + \hat{\alpha}_2 + residual_2 = 0$$

$$(16.1.11) \quad \hat{\beta}_{31} output + \hat{\beta}_{32} labor - capital + \hat{\alpha}_3 + residual_3 = 0$$

*This gives the following table of coefficients:*

	output	labor	capital	intercept
	-1	0.8072782	0.2330535	-0.1773097
	0.73812541	-1	-0.01105754	1.27424214
	1.67579357	-0.08695944	-1	-2.72052726

*Now divide the second and third rows by the negative of their first coefficient, so that the coefficient of output becomes -1:*

out	labor	capital	intercept
-1	0.8072782	0.2330535	-0.1773097
-1	1/0.73812541	0.01105754/0.73812541	-1.27424214/0.73812541
-1	0.08695944/1.67579357	1/1.67579357	2.72052726/1.67579357

*After performing the divisions the following numbers are obtained:*

output	labor	capital	intercept
-1	0.8072782	0.2330535	-0.1773097
-1	1.3547833	0.014980570	-1.726322
-1	0.05189149	0.59673221	1.6234262

*These results can also be re-written in the form given by Table 2.*

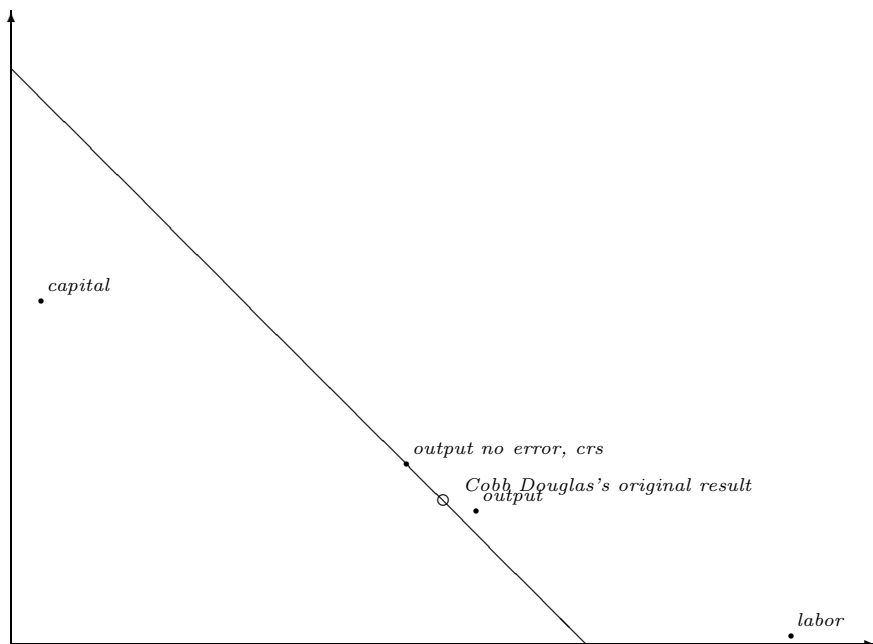


FIGURE 2. Coefficients of capital (vertical) and labor (horizontal), dependent variable output, 1899–1922

	Intercept	Slope of output wrt labor	Slope of output wrt capital
Regression of output on labor and capital			
Regression of labor on output and capital			
Regression of capital on output and labor			

TABLE 2. Comparison of coefficients in elementary regressions

Fill in the values for the whole period and also for several sample subperiods. Make a scatter plot of the contents of this table, i.e., represent each regression result as a point in a plane, using different colors for different sample periods.

PROBLEM 222. Given a univariate problem with three variables all of which have zero mean, and a linear constraint that the coefficients of all variables sum to 0. (This

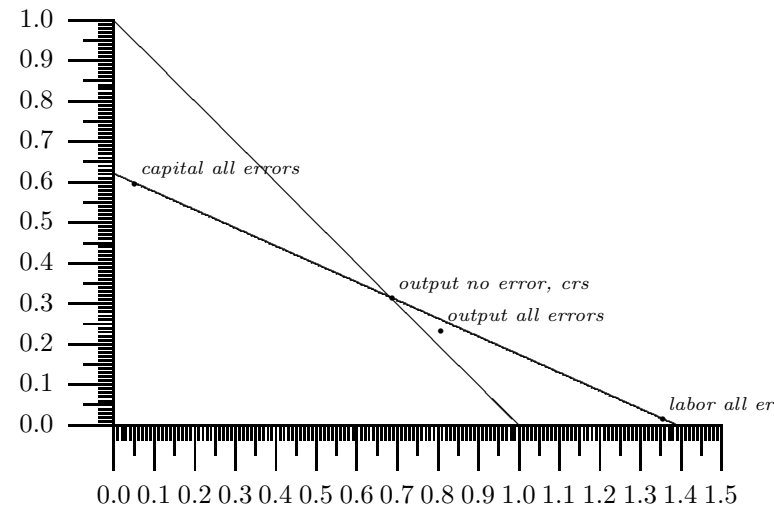


FIGURE 3. Coefficient of capital (vertical) and labor (horizontal) in the elementary regressions, dependent variable output, 1899–1922

is the model apparently appropriate to the Cobb-Douglas data, with the assumption of constant returns to scale, after taking out the means.) Call the observed variables  $x$ ,  $y$ , and  $z$ , with underlying systematic variables  $x^*$ ,  $y^*$ , and  $z^*$ , and errors  $u$ ,  $v$ , and  $w$ .

- a. Write this model in the form (??).

ANSWER.

$$(16.1.12) \quad \begin{bmatrix} x^* & y^* & z^* \end{bmatrix} \begin{bmatrix} -1 \\ \beta \\ 1 - \beta \end{bmatrix} = 0 \quad \text{or} \quad \begin{aligned} x^* &= \beta y^* + (1 - \beta)z^* \\ x &= x^* + u \\ y &= y^* + v \\ z &= z^* + w. \end{aligned}$$

- b. The moment matrix of the systematic variables can be written fully in terms of  $\sigma_{y^*}^2$ ,  $\sigma_{z^*}^2$ ,  $\sigma_{y^*z^*}$ , and the unknown parameter  $\beta$ . Write out the moment matrix and therefore the Frisch decomposition.

ANSWER. The moment matrix is the middle matrix in the following Frisch decomposition:

$$(16.1.13) \quad \begin{bmatrix} \sigma_x^2 & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_y^2 & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_z^2 \end{bmatrix} =$$



$$(16.1.14) \quad = \begin{bmatrix} \beta^2 \sigma_{y^*}^2 + 2\beta(1-\beta)\sigma_{y^*z^*} + (1-\beta)^2 \sigma_{z^*}^2 & \beta \sigma_{y^*}^2 + (1-\beta)\sigma_{y^*z^*} & \beta \sigma_{y^*z^*} + (1-\beta)\sigma_{z^*}^2 \\ \beta \sigma_{y^*}^2 + (1-\beta)\sigma_{y^*z^*} & \sigma_{y^*}^2 & \sigma_{y^*z^*} \\ \beta \sigma_{y^*z^*} + (1-\beta)\sigma_{z^*}^2 & \sigma_{y^*z^*} & \sigma_{z^*}^2 \end{bmatrix} + \begin{bmatrix} \sigma_u^2 & 0 \\ 0 & \sigma_v^2 \\ 0 & 0 \end{bmatrix} \quad \square$$

• c. Show that the unknown parameters are not yet identified. However, if one makes the additional assumption that one of the three error variances  $\sigma_u^2$ ,  $\sigma_v^2$ , or  $\sigma_w^2$  is zero, then the equations are identified. Since the quantity of output presumably has less error than the other two variables, assume  $\sigma_u^2 = 0$ . Under this assumption, show that

$$(16.1.15) \quad \beta = \frac{\sigma_x^2 - \sigma_{xz}}{\sigma_{xy} - \sigma_{xz}}$$

and this can be estimated by replacing the variances and covariances by their sample counterparts. In a similar way, derive estimates of all other parameters of the model.

ANSWER. Solving (16.1.14) one gets from the yz element of the covariance matrix

$$(16.1.16) \quad \sigma_{y^*z^*} = \sigma_{yz}$$

and from the xz element

$$(16.1.17) \quad \sigma_{z^*}^2 = \frac{\sigma_{xz} - \beta \sigma_{yz}}{1 - \beta}$$

Similarly, one gets from the xy element:

$$(16.1.18) \quad \sigma_{y^*}^2 = \frac{\sigma_{xy} - (1 - \beta)\sigma_{yz}}{\beta}$$

Now plug (16.1.16), (16.1.17), and (16.1.18) into the equation for the xx element:

$$(16.1.19) \quad \sigma_x^2 = \beta(\sigma_{xy} - (1 - \beta)\sigma_{yz}) + 2\beta(1 - \beta)\sigma_{yz} + (1 - \beta)(\sigma_{xz} - \beta\sigma_{yz}) + \sigma_u^2$$

$$(16.1.20) \quad = \beta\sigma_{xy} + (1 - \beta)\sigma_{xz} + \sigma_u^2$$

Since we are assuming  $\sigma_u^2 = 0$  this last equation can be solved for  $\beta$ :

$$(16.1.21) \quad \beta = \frac{\sigma_x^2 - \sigma_{xz}}{\sigma_{xy} - \sigma_{xz}}$$

If we replace the variances and covariances by the sample variances and covariances, this gives an estimate of  $\beta$ . □

• d. Evaluate these formulas numerically. In order to get the sample means and the sample covariance matrix of the data, you may issue the SAS commands

```
proc corr cov nocorr data=cdtmp;
var logout loglab logcap;
run;
```

These commands are in the file *cbbcovma.sas* on the disk.

ANSWER. Mean vector and covariance matrix are

$$(16.1.22) \quad \begin{bmatrix} \text{LOGOUT} \\ \text{LOGLAB} \\ \text{LOGCAP} \end{bmatrix} \sim \left( \begin{bmatrix} 5.07734 \\ 4.96272 \\ 5.35648 \end{bmatrix}, \begin{bmatrix} 0.0724870714 & 0.0522115563 & 0.1169330807 \\ 0.0522115563 & 0.0404318579 & 0.0839798588 \\ 0.1169330807 & 0.0839798588 & 0.2108441826 \end{bmatrix} \right)$$

Therefore equation (16.1.15) gives

$$(16.1.23) \quad \hat{\beta} = \frac{0.0724870714 - 0.1169330807}{0.0522115563 - 0.1169330807} = 0.686726861149148$$

In Figure 3, the point  $(\hat{\beta}, 1 - \hat{\beta})$  is exactly the intersection of the long dotted line with the constraint line.

• e. The fact that all 3 points lie almost on the same line indicates that there may be 2 linear relations: log labor is a certain coefficient times log output, and log capital is a different coefficient times log output. I.e.,  $y^* = \delta_1 + \gamma_1 x^*$  and  $z^* = \delta_2 + \gamma_2 x^*$ . In other words, there is no substitution. What would be the two coefficients  $\gamma_1$  and  $\gamma_2$  if this were the case?

ANSWER. Now the Frisch decomposition is

$$(16.1.24) \quad \begin{bmatrix} \sigma_x^2 & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_y^2 & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_z^2 \end{bmatrix} = \sigma_{x^*}^2 \begin{bmatrix} 1 & \gamma_1 & \gamma_2 \\ \gamma_1 & \gamma_1^2 & \gamma_1 \gamma_2 \\ \gamma_2 & \gamma_1 \gamma_2 & \gamma_2^2 \end{bmatrix} + \begin{bmatrix} \sigma_u^2 & 0 & 0 \\ 0 & \sigma_v^2 & 0 \\ 0 & 0 & \sigma_w^2 \end{bmatrix}.$$

Solving this gives (obtain  $\gamma_1$  by dividing the 32-element by the 31-element,  $\gamma_2$  by dividing 32-element by the 12-element,  $\sigma_{x^*}^2$  by dividing the 21-element by  $\gamma_1$ , etc.)

$$(16.1.25) \quad \begin{aligned} \gamma_1 &= \frac{\sigma_{yz}}{\sigma_{xy}} = \frac{0.0839798588}{0.1169330807} = 0.7181873452513939 & \sigma_u^2 &= \sigma_x^2 - \frac{\sigma_{yx}\sigma_{xz}}{\sigma_{yz}} = 0.0724870714 \\ \gamma_2 &= \frac{\sigma_{yz}}{\sigma_{xz}} = \frac{0.0839798588}{0.0522115563} = 1.608453467992104 & \sigma_v^2 &= \sigma_y^2 - \frac{\sigma_{xy}\sigma_{yz}}{\sigma_{xz}} \\ \sigma_{x^*}^2 &= \frac{\sigma_{yx}\sigma_{xz}}{\sigma_{yz}} = \frac{0.0522115563 \cdot 0.1169330807}{0.0839798588} = 0.0726990758 & \sigma_w^2 &= \sigma_z^2 - \frac{\sigma_{xz}\sigma_{zy}}{\sigma_{xy}} \end{aligned}$$

This model is just barely rejected by the data since it leads to a slightly negative variance for  $\mathbf{U}$ .

• f. The assumption that there are two linear relations is represented as light-blue line in Figure 3. What is the equation of this line?

ANSWER. If  $y = \gamma_1 x$  and  $z = \gamma_2 x$  then the equation  $x = \beta_1 y + \beta_2 z$  holds whenever  $\beta_1 \gamma_1 + \beta_2 \gamma_2 = 1$ . This is a straight line in the  $\beta_1, \beta_2$ -plane, going through the points and  $(0, 1/\gamma_2)$ ,  $(0, \frac{0.0522115563}{0.0839798588} = 0.6217152189353289)$  and  $(1/\gamma_1, 0) = (\frac{0.1169330807}{0.0839798588} = 1.3923943475361023, 0)$ . This line is in the figure, and it is just a tiny bit on the wrong side of the dotted line connecting the two estimates.

## 16.2. Houthakker's Data

For this example we will use Berndt's textbook [Ber91], which discusses some of the classic studies in the econometric literature.

One example described there is the estimation of a demand function for electricity [Hou51], which is the first multiple regression with several variables run on a computer. In this exercise you are asked to do all steps in exercise 1 and 3 in chapter 7 of Berndt, and use the additional facilities of R to perform other steps of data analysis which Berndt did not ask for, such as, for instance, explore the best subset of regressors using `leaps` and the best nonlinear transformation using `avas`, do some diagnostics, search for outliers or influential observations, and check the normality of residuals by a probability plot.

**PROBLEM 223.** *4 points* The electricity demand data from [Hou51] are available on the web at `www.econ.utah.edu/ehrbar/data/ukelec.txt`. Import these data into your favorite statistics package. For R you need the command `ukelec <- read.table("http://www.econ.utah.edu/ehrbar/data/ukelec.txt")`. Make a scatterplot matrix of these data using e.g. `pairs(ukelec)` and describe what you see.

**ANSWER.** `inc` and `cap` are negatively correlated. `cap` is capacity of rented equipment and not equipment owned. Apparently customers with higher income buy their equipment instead of renting it.

`gas6` and `gas8` are very highly correlated. `mc4`, `mc6`, and `mc8` are less highly correlated, the correlation between `mc6` and `mc8` is higher than that between `mc4` and `mc6`. It seems electricity prices have been coming down.

`kwh`, `inc`, and `exp` are strongly positively correlated.

The stripes in all the plots which have `mc4`, `mc6`, or `mc8` in them come from the fact that the marginal cost of electricity is a round number.

Electricity prices and `kwh` are negatively correlated.

There is no obvious positive correlation between `kwh` and `cap` or `exp` and `cap`.

Prices of electricity and gas are somewhat positively correlated, but not much.

When looking at the correlations of `inc` with the other variables, there are several outliers which could have a strong "leverage" effect.

In 1934, those with high income had lower electricity prices than those with low income. This effect dissipated by 1938.

No strong negative correlations anywhere.

`cust` negatively correlated with `inc`, because rich people live in smaller cities? □

If you simply type `ukelec` in R, it will print the data on the screen. The variables have the following meanings:

`cust` Average number of consumers with two-part tariffs for electricity in 1937–38, in thousands. Two-part tariff means: they pay a fixed monthly sum plus a certain "running charge" times the number of kilowatt hours they use.

`inc` Average income of two-part consumers, in pounds per year. (Note that one pound had 240 pence at that time.)

`mc4` The running charge (marginal cost) on domestic two-part tariffs in 1933–38, in pence per KWH. (The marginal costs are the costs that depend on the number of kilowatt hours only, it is the cost of one additional kilowatt hour.)

`mc6` The running charge (marginal cost) on domestic two-part tariffs in 1935–38, in pence per KWH

`mc8` The running charge (marginal cost) on domestic two-part tariffs in 1937–38, in pence per KWH

`gas6` The marginal price of gas in 1935–36, in pence per therm

`gas8` The marginal price of gas in 1937–38, in pence per therm

`kwh` Consumption on domestic two-part tariffs per consumer in 1937–38, in kilowatt hours

`cap` The average holdings (capacity) of heavy electric equipment bought on hire or purchase (leased) by domestic two-part consumers in 1937–38, in kilowatts

`expen` The average total expenditure on electricity by two-part consumers in 1937–38, in pounds

The function `summary(ukelec)` displays summary statistics about every variable in the data frame.

Since every data frame in R is a list, it is possible to access the variables in `ukelec` by typing `ukelec$mc4` etc. Try this; if you type this and then a return, you will get a listing of `mc4`. In order to have all variables available as separate objects and save space by not typing `ukelec$` all the time, one has to "mount" the data frame by the command `attach(ukelec)`. After this, the individual data series can simply be printed on the screen by typing the name of the variable, for instance `mc4`, and then the return key.

**PROBLEM 224.** *2 points* Make boxplots of `mc4`, `mc6`, and `mc8` in the same graphics window next to each other, and the same with `gas6` and `gas8`.

**PROBLEM 225.** *2 points* How would you answer the question whether marginal gas prices vary more or less than those of electricity (say in the year 1936)?

**ANSWER.** Marginal gas prices vary a little more than electricity prices, although electricity prices were the newer technology, and although gas prices are much more stable over time than the electricity prices. Compare `sqrt(var(mc6))/mean(mc6)` with `sqrt(var(gas6))/mean(gas6)`. You get 0.176 versus 0.203. Another way would be to compute `max(mc6)/min(mc6)` and compare with `max(gas6)/min(gas6)`: you get 2.27 versus 2.62. In any case this is a lot of variation.

**PROBLEM 226.** *2 points* Make a plot of the (empirical) density function of `inc` and `gas6` and interpret the results.

**PROBLEM 227.** *2 points* Is electricity a big share of total income? Which command is better: `mean(expen/inc)` or `mean(expen)/mean(inc)`? What other options

are there? Actually, there is a command which is clearly better than at least one of the above, can you figure out what it is?

ANSWER. The proportion is small, less than 1 percent. The two above commands give 0.89% and 0.84%. The command `sum(cust*expen) / sum(cust*inc)` is better than `mean(expen) / mean(inc)`, because each component in `expen` and `inc` is the mean over many households, the number of households given by `cust`. `mean(expen)` is therefore an average over averages over different population sizes, not a good idea. `sum(cust*expen)` is total expenditure in all households involved, and `sum(cust*inc)` is total income in all households involved. `sum(cust*expen) / sum(cust*inc)` gives the value 0.92%. Another option is `median(expen/inc)` which gives 0.91%. A good way to answer this question is to plot it: `plot(expen,inc)`. You get the line where expenditure is 1 percent of income by `abline(0,0.01)`. For higher incomes expenditure for electricity levels off and becomes a lower share of income. □

PROBLEM 228. Have your computer compute the sample correlation matrix of the data. The R-command is `cor(ukelec)`

- a. 4 points Are there surprises if one looks at the correlation matrix?

ANSWER. Electricity consumption `kwh` is slightly negatively correlated with gas prices and with the capacity. If one takes the correlation matrix of the logarithmic data, one gets the expected positive signs.

marginal prices of gas and electricity are positively correlated in the order of 0.3 to 0.45.

higher correlation between `mc6` and `mc8` than between `mc4` and `mc6`.

Correlation between `expen` and `cap` is negative and low in both matrices, while one should expect positive correlation. But in the logarithmic matrix, `mc6` has negative correlation with `expen`, i.e., elasticity of electricity demand is less than 1.

In the logarithmic data, `cust` has higher correlations than in the non-logarithmic data, and it is also more nearly normally distributed.

`inc` has negative correlation with `mc4` but positive correlation with `mc6` and `mc8`. (If one looks at the scatterplot matrix this seems just random variations in an essentially zero correlation).

`mc6` and `expen` are positively correlated, and so are `mc8` and `expen`. This is due to the one outlier with high `expen` and high income and also high electricity prices.

The marginal prices of electricity are not strongly correlated with `expen`, and in 1934, they are negatively correlated with `income`.

From the scatter plot of `kwh` versus `cap` it seems there are two datapoints whose removal might turn the sign around. To find out which they are do `plot(kwh,cap)` and then use the identify function: `identify(kwh,cap,labels=row.names(ukelec))`. The two outlying datapoints are Halifax and Wallase. Wallase has the highest income of all towns, namely, 1422, while Halifax's income of 352 is close to the minimum, which is 279. High income customers do not lease their equipment but buy it. □

- b. 3 points The correlation matrix says that `kwh` is negatively related with `cap`, but the correlation of the logarithm gives the expected positive sign. Can you explain this behavior?

ANSWER. If one plots the data using `plot(cap,kwh)` one sees that the negative correlation comes from the two outliers. In a logarithmic scale, these two are no longer so strong outliers. □

PROBLEM 229. Berndt on p. 338 defines the intramarginal expenditure  $f$  as  $\text{expen} - \text{mc8} * \text{kwh} / 240$ . What is this, and what do you find out looking at it?

After this preliminary look at the data, let us run the regressions.

PROBLEM 230. 6 points Write up the main results from the regressions which R are run by the commands

```
houth.olsfit <- lm(formula = kwh ~ inc+I(1/mc6)+gas6+cap)
houth.glsfit <- lm(kwh ~ inc+I(1/mc6)+gas6+cap, weight=cust)
houth.olsloglogfit <- lm(log(kwh) ~
log(inc)+log(mc6)+log(gas6)+log(cap))
```

Instead of `1/mc6` you had to type `I(1/mc6)` because the slash has a special meaning in formulas, creating a nested design, therefore it had to be “protected” by applying the function `I()` to it.

If you then type `houth.olsfit`, a short summary of the regression results will be displayed on the screen. There is also the command `summary(houth.olsfit)`, which gives you a more detailed summary. If you type `plot(houth.olsfit)` you will get a series of graphics relevant for this regression.

ANSWER. All the expected signs.

Gas prices do not play a great role in determining electricity consumption, despite the “corner” Berndt talks about on p. 337. Especially the logarithmic regression makes gas prices highly insignificant!

The weighted estimation has a higher  $R^2$ .

PROBLEM 231. 2 points The output of the OLS regression gives as standard error of `inc` the value of 0.18, while in the GLS regression it is 0.20. For the other variables, the standard error as given in the GLS regression is lower than that in the OLS regression. Does this mean that one should use for `inc` the OLS estimate and for the other variables the GLS estimates?

PROBLEM 232. 5 points Show, using the `leaps` procedure on R or some other selection of regressors, that the variables Houthakker used in his GLS-regression are the “best” among the following: `inc`, `mc4`, `mc6`, `mc8`, `gas6`, `gas8`, `cap` using either the  $C_p$  statistic or the adjusted  $R^2$ . (At this stage, do not transform the variables but just enter them into the regression untransformed, but do use the weights, which are theoretically well justified).

To download the `leaps` package, use `install.packages("leaps", lib="C:/Documents and Settings/420lab.420LAB/My Documents")` and to call it up, use `library(leaps, lib.loc="C:/Documents and Settings/420lab.420LAB/My Documents")`. If the library `ecmet` is available, the command `ecmet.script(houthsel)` runs the following script:

```

library(leaps)
data(ukelec)
attach(ukelec)
houth.glsleaps<-leaps(x=cbind(inc,mc4,mc6,mc8,gas6,gas8,cap),
                    y=kwh, wt=cust, method="Cp",
                    nbest=5, strictly.compatible=F)
ecmet.prompt("Plot Mallow's Cp against number of regressors:")
plot(houth.glsleaps$size, houth.glsleaps$Cp)
ecmet.prompt("Throw out all regressions with a Cp > 50 (big gap)")
plot(houth.glsleaps$size[houth.glsleaps$Cp<50],
     houth.glsleaps$Cp[houth.glsleaps$Cp<50])
ecmet.prompt("Cp should be roughly equal the number of regressors")
abline(0,1)
cat("Does this mean the best regression is overfitted?")
ecmet.prompt("Click at the points to identify them, left click to quit")
## First construct the labels
lngth <- dim(houth.glsleaps$which)[1]
included <- as.list(1:lngth)
for (ii in 1:lngth)
  included[[ii]] <- paste(
    colnames(houth.glsleaps$which)[houth.glsleaps$which[ii,]],
    collapse=",")
identify(x=houth.glsleaps$size, y=houth.glsleaps$Cp, labels=included)
ecmet.prompt("Now use regsubsets instead of leaps")
houth.glsrss<- regsubsets.default(x=cbind(inc,mc4,mc6,mc8,gas6,gas8,cap),
                               y=kwh, weights=cust, method="exhaustive")
print(summary.regsubsets(houth.glsrss))
plot.regsubsets(houth.glsrss, scale="Cp")
ecmet.prompt("Now order the variables")
houth.glsrsord<- regsubsets.default(x=cbind(inc,mc6,cap,gas6,gas8,mc8,mc4),
                                   y=kwh, weights=cust, method="exhaustive")
print(summary.regsubsets(houth.glsrsord))
plot.regsubsets(houth.glsrsord, scale="Cp")

```

PROBLEM 233. Use *avas* to determine the “best” nonlinear transformations of the explanatory and the response variable. Since the weights are theoretically well justified, one should do it for the weighted regression. Which functions do you think one should use for the different regressors?

PROBLEM 234. 3 points Then, as a check whether the transformation interfered with data selection, redo *leaps*, but now with the transformed variables. Show that

the GLS-regression Houthakker actually ran is the “best” regression among the following variables: *inc*,  $1/mc4$ ,  $1/mc6$ ,  $1/mc8$ , *gas6*, *gas8*, *cap* using either  $C_p$  statistic or the adjusted  $R^2$ .

PROBLEM 235. Diagnostics, the identification of outliers or influential observations is something which we can do easily with R, although Berndt did not ask for it. The command `houth.glsinf<-lm.influence(houth.glsfit)` gives you the building blocks for many of the regression diagnostics statistics. Its output is a list of objects: A matrix whose rows are all the the least squares estimates  $\hat{\beta}(i)$  when the  $i$ th observation is dropped, a vector with all the  $s(i)$ , and a vector with all the  $h(i)$ . A more extensive function is `influence.measures(houth.glsfit)`, it has Cook's distance and others.

In order to look at the residuals, use the command `plot(resid(houth.glsfit), type="h")` or `plot(rstandard(houth.glsfit), type="h")` or `plot(rstudent(houth.glsfit), type="h")`. To add the axis do `abline(0,0)`. If you wanted to check the residuals for normality, you would use `qqnorm(rstandard(houth.glsfit))`.

PROBLEM 236. Which commands do you need to plot the predictive residuals?

PROBLEM 237. 4 points Although there is good theoretical justification for using *cust* as weights, one might wonder if the data bear this out. How can you check this?

ANSWER. Do `plot(cust, rstandard(houth.olsfit))` and `plot(cust, rstandard(houth.glsfit))`. In the first plot, smaller numbers of customers have larger residuals, in the second plot this is mitigated. Also the OLS plot has two terrible outliers, which are brought more into range with GLS.

PROBLEM 238. The variable *cap* does not measure the capacity of all electrical equipment owned by the households, but only those appliances which were leased from the Electric Utility company. A plot shows that people with higher income do not lease as much but presumably purchase their appliances outright. Does this mean that *cap* variable should not be in the regression?

### 16.3. Long Term Data about US Economy

The dataset `us1t` is described in [DL91]. Home page of the authors is `www.ces.iupui.edu/us1t` has the variables *kn*, *kg* (net and gross capital stock in current \$), *kn2*, *kg2* (the same in 1982\$), *hours* (hours worked), *wage* (hourly wage in current dollars), *gnp*, *gnp2*, *nnp*, *inv2* (investment in 1982 dollars), *r* (profit rate ( $nnp - wage \times hours$ )/*kn*), *u* (capacity utilization), *kne*, *kge*, *kne2*, *kge2*, *inve2* (capital stock and investment data for equipment), *kns*, *kgs*, *kns2*, *kgs2*, *invs2* (the same for structures).

Capital stock data were estimated separately for structures and equipment and then added up, i.e.,  $kn2 = kne2 + kns2$  etc. Capital stock since 1925 has been constructed from annual investment data, and prior to 1925 the authors of the series

apparently went the other direction: they took someone's published capital stock estimates and constructed investment from it. In the 1800s, only a few observations were available, which were then interpolated. The capacity utilization ratio is equal to the ratio of `gnp2` to its trend, i.e., it may be negative.

Here are some possible commands for your R-session: `data(uslt)` makes the data available; `uslt.clean<-na.omit(uslt)` removes missing values; this dataset starts in 1869 (instead of 1805). `attach(uslt.clean)` makes the variables in this dataset available. Now you can plot various series, for instance `plot((nnp-hours*wage)/nnp, type="l")` plots the profit share, or `plot(gnp/gnp2, kg/kg2, type="l")` gives you a scatter plot of the price level for capital goods versus that for gnp. The command `plot(r, kn2/hours, type="b")` gives both points and dots; `type = "o"` will have the dots overlaid the line. After the plot you may issue the command `identify(r, kn2/hours, label=1869:1989)` and then click with the left mouse button on the plot those data points for which you want to have the years printed.

If you want more than one timeseries on the same plot, you may do `matplot(1869:1989, cbind(kn2,kns2), type="l")`. If you want the y-axis logarithmic, say `matplot(1869:1989, cbind(gnp/gnp2,kns/kns2,kne/kne2), type="l", log="y")`.

**PROBLEM 239.** *Computer assignment: Make a number of such plots on the screen, and import the most interesting ones into your wordprocessor. Each class participant should write a short paper which shows the three most interesting plots, together with a written explanation why these plots seem interesting.*

To use `pairs` or `xgobi`, you should carefully select the variables you want to include, and then you need the following preparations: `usltsplom <- cbind(gnp2=gnp2, kn2=kn2, inv2=inv2, hours=hours, year=1869:1989) dimnames(usltsplom)[[1]] <- paste(1869:1989)` The `dimnames` function adds the row labels to the matrix, so that you can see which year it is. `pairs(usltsplom)` or `library(xgobi)` and then `xgobi(usltsplom)`

You can also run regressions with commands of the following sort: `lm.fit <- lm(formula = gnp2 ~ hours + kne2 + kns2)`. You can also fit a “generalized additive model” with the formula `gam.fit <- gam(formula = gnp2 ~ s(hours) + s(kne2) + s(kns2))`. This is related to the `avas` command we talked about in class. It is discussed in [CH93].

## 16.4. Dougherty Data

We have a new dataset, in both SAS and Splus, namely the data described in [Dou92].

There are more data than in the tables at the end of the book; `prelcosm` for instance is the relative price of cosmetics, it is `100*pcosm/ptpe`, but apparently truncated at 5 digits.

## 16.5. Wage Data

The two datasets used in [Ber91, pp. 191–209] are available in R as the data frames `cps78` and `cps85`. In R on unix, the data can be downloaded by `cps78 <- readSDML("http://www.econ.utah.edu/ehrbar/data/cps78.sdml")`, and `cps85 <- readSDML("http://www.econ.utah.edu/ehrbar/data/cps85.sdml")`, and `attach(cps78)` corresponding for `cps85`. The original data provided by Berndt contain many dummy variables. The data frames in R have the same data coded as “factor” variables instead of dummies. These “factor” variables automatically generate dummies which are included in the `model` statement.

`cps78` consists of 550 randomly selected employed workers from the May 1978 current population survey, and `cps85` consists of 534 randomly selected employed workers from the May 1985 current population survey. These are surveys of 50,000 households conducted monthly by the U.S. Department of Commerce. They serve as the basis for the national employment and unemployment statistics. Data are collected on a number of individual characteristics as well as employment statistics. The present extracts were performed by Leslie Sundt of the University of Arizona.

`ed` = years of education  
`ex` = years of labor market experience (= `age - ed - 6`, or 0 if this is a negative number).

`lnwage` = natural logarithm of average hourly earnings  
`age` = age in years  
`ndep` = number of dependent children under 18 in household (only in `cps78`)  
`region` has levels North, South  
`race` has levels Other, Nonwhite, Hispanic. Nonwhite is mainly the Blacks, and Other is mainly the Non-Hispanic Whites.  
`gender` has levels Male, Female  
`marr` has levels Single, Married  
`union` has levels Nonunion, Union  
`industry` has levels Other, Manuf, and Constr  
`occupation` has levels Other, Manag, Sales, Cler, Serv, and Prof  
 Here is a log of my commands for exercises 1 and 2 in [Ber91, pp. 194–197].

```
> cps78 <- readSDML("http://www.econ.utah.edu/ehrbar/data/cps78.sdml")
> attach(cps78)
> ###Exercise 1a (2 points) in chapter V of Berndt, p. 194
> #Here is the arithmetic mean of hourly wages:
> mean(exp(lnwage))
[1] 6.062766
> #Here is the geometric mean of hourly wages:
> #(Berndt's instructions are apparently mis-formulated):
> exp(mean(lnwage))
```

```
[1] 5.370935
> #Geometric mean is lower than arithmetic, due to Jensen's inequality
> #if the year has 2000 hours, this gives an annual wage of
> 2000*exp(mean(lnwage))
[1] 10741.87
> #What are arithmetic mean and standard deviation of years of schooling
> #and years of potential experience?
> mean(ed)
[1] 12.53636
> sqrt(var(ed))
[1] 2.772087
> mean(ex)
[1] 18.71818
> sqrt(var(ex))
[1] 13.34653
> #experience has much higher standard deviation than education, not surpr
> ##Exercise 1b (1 point) can be answered with the two commands
> table(race)
  Hisp Nonwh Other
    36   57  457
> table(race, gender)
      gender
race  Female Male
  Hisp     12   24
  Nonwh     28   29
  Other    167  290
> #Berndt also asked for the sample means of certain dummy variables;
> #This has no interest in its own right but was an intermediate
> #step in order to compute the numbers of cases as above.
> ##Exercise 1c (2 points) can be answered using tapply
> tapply(ed,gender,mean)
  Female   Male
12.76329 12.39942
> #now the standard deviation:
> sqrt(tapply(ed,gender,var))
  Female   Male
2.220165 3.052312
> #Women do not have less education than men; it is about equal,
> #but their standard deviation is smaller
> #Now the geometric mean of the wage rate:
```

```
> exp(tapply(lnwage,gender,mean))
  Female   Male
4.316358 6.128320
> #Now do the same with race
> ##Exercise 1d (4 points)
> detach()
> ##This used to be my old command:
> cps85 <- read.table("~/dpkg/ecmet/usr/share/ecmet/usr/lib/R/librari
> #But this should work for everyone (perhaps only on linux):
> cps85 <- readSDML("http://www.econ.utah.edu/ehrbart/data/cps85.sdm
> attach(cps85)
> mean(exp(lnwage))
[1] 9.023947
> sqrt(var(lnwage))
[1] 0.5277335
> exp(mean(lnwage))
[1] 7.83955
> 2000*exp(mean(lnwage))
[1] 15679.1
> 2000*exp(mean(lnwage))/1.649
[1] 9508.248
> #real wage has fallen
> tapply(exp(lnwage), gender, mean)
  Female   Male
7.878743 9.994794
> tapply(exp(lnwage), gender, mean)/1.649
  Female   Male
4.777891 6.061125
> #Compare that with 4.791237 6.830132 in 1979:
> #Male real wages dropped much more than female wages
> ##Exercise 1e (3 points)
> #using cps85
> w <- mean(lnwage); w
[1] 2.059181
> s <- sqrt(var(lnwage)); s
[1] 0.5277335
> lnwagef <- factor(cut(lnwage, breaks = w+s*c(-4,-2,-1,0,1,2,4)))
> table(lnwagef)
lnwagef
(-0.0518,1] (1,1.53] (1.53,2.06] (2.06,2.59] (2.59,3.11] (3.11,4
```

```

      3      93      174      180      72      12
> ks.test(lnwage, "pnorm")

One-sample Kolmogorov-Smirnov test

data: lnwage
D = 0.8754, p-value = < 2.2e-16
alternative hypothesis: two.sided

> ks.test(lnwage, "pnorm", mean=w, sd =s)

One-sample Kolmogorov-Smirnov test

data: lnwage
D = 0.0426, p-value = 0.2879
alternative hypothesis: two.sided

> #Normal distribution not rejected
>
> #If we do the same thing with
> wage <- exp(lnwage)
> ks.test(wage, "pnorm", mean=mean(wage), sd =sqrt(var(wage)))

One-sample Kolmogorov-Smirnov test

data: wage
D = 0.1235, p-value = 1.668e-07
alternative hypothesis: two.sided

> #This is highly significant, therefore normality rejected
>
> #An alternative, simpler way to answer question 1e is by using qqnorm
> qqnorm(lnwage)
> qqnorm(exp(lnwage))
> #Note that the SAS proc univariate rejects that wage is normally distri
> #but does not reject that lnwage is normally distributed.
> ###Exercise 2a (3 points), p. 196
> summary(lm(lnwage ~ ed, data = cps78))

```

```

Call:
lm(formula = lnwage ~ ed, data = cps78)

Residuals:
    Min       1Q   Median       3Q      Max
-2.123168 -0.331368 -0.007296  0.319713  1.594445

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.030445   0.092704  11.115 < 2e-16 ***
ed           0.051894   0.007221   7.187 2.18e-12 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.469 on 548 degrees of freedom
Multiple R-Squared:  0.08613, Adjusted R-squared:  0.08447
F-statistic: 51.65 on 1 and 548 degrees of freedom, p-value: 2.181e-12

> #One year of education increases wages by 5 percent, but low R^2.
> #Mincer (5.18) had 7 percent for 1959
> #Now we need a 95 percent confidence interval for this coefficient
> 0.051894 + 0.007221*qt(0.975, 548)
[1] 0.06607823
> 0.051894 - 0.007221*qt(0.975, 548)
[1] 0.03770977
> ##Exercise 2b (3 points): Include union participation
> summary(lm(lnwage ~ union + ed, data=cps78))

Call:
lm(formula = lnwage ~ union + ed, data = cps78)

Residuals:
    Min       1Q   Median       3Q      Max
-2.331754 -0.294114  0.001475  0.263843  1.678532

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.859166   0.091630   9.376 < 2e-16 ***
unionUnion  0.305129   0.041800   7.300 1.02e-12 ***
ed           0.058122   0.006952   8.361 4.44e-16 ***

```

```

---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.4481 on 547 degrees of freedom
Multiple R-Squared:  0.1673, Adjusted R-squared:  0.1642
F-statistic: 54.93 on 2 and 547 degrees of freedom, p-value:    0

> exp(0.058)
[1] 1.059715
> exp(0.305129)
[1] 1.3568
> # Union members have 36 percent higher wages
> # The test whether union and nonunion members have the same intercept
> # is the same as the test whether the union dummy is 0.
> # t-value = 7.300 which is highly significant,
> # i.e., they are different.

> #The union variable is labeled unionUnion, because
> #it is labeled 1 for Union and 0 for Nonun. Check with the command
> contrasts(cps78$union)
      Union
Nonun      0
Union      1
> #One sees it also if one runs
> model.matrix(lnwage ~ union + ed, data=cps78)
      (Intercept) union ed
1             1      0 12
2             1      1 12
3             1      1  6
4             1      1 12
5             1      0 12
> #etc, rest of output flushed
> #and compares this with
> cps78$union[1:5]
[1] Nonun Union Union Union Nonun
Levels:  Nonun Union
> #Consequently, the intercept for nonunion is 0.8592
> #and the intercept for union is 0.8592+0.3051=1.1643.
> #Can I have a different set of dummies constructed from this factor?
> #We will first do

```

```

> ##Exercise 2e (2 points)
> contrasts(union)<-matrix(c(1,0),nrow=2,ncol=1)
> #This generates a new contrast matrix
> #which covers up that in cps78
> #Note that I do not say "data=cps78" in the next command:
> summary(lm(lnwage ~ union + ed))

Call:
lm(formula = lnwage ~ union + ed)

Residuals:
      Min       1Q   Median       3Q      Max
-2.331754 -0.294114  0.001475  0.263843  1.678532

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.164295   0.090453  12.872 < 2e-16 ***
union1       -0.305129   0.041800  -7.300 1.02e-12 ***
ed           0.058122   0.006952   8.361 4.44e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.4481 on 547 degrees of freedom
Multiple R-Squared:  0.1673, Adjusted R-squared:  0.1642
F-statistic: 54.93 on 2 and 547 degrees of freedom, p-value:    0

> #Here the coefficients are different,
> #but it is consistent with the above result.
> ##Exercise 2c (2 points):  If I want to have two contrasts from the
> contrasts(union,2)<-matrix(c(1,0,0,1),nrow=2,ncol=2)
> #The additional argument 2
> #specifies different number of contrasts than it expects
> #Now I have to suppress the intercept in the regression
> summary(lm(lnwage ~ union + ed - 1))

Call:
lm(formula = lnwage ~ union + ed - 1)

Residuals:
      Min       1Q   Median       3Q      Max

```



```
-2.331754 -0.294114 0.001475 0.263843 1.678532
```

```
Coefficients:
```

```
      Estimate Std. Error t value Pr(>|t|)
union1 0.859166  0.091630  9.376 < 2e-16 ***
union2 1.164295  0.090453 12.872 < 2e-16 ***
ed      0.058122  0.006952  8.361 4.44e-16 ***
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.4481 on 547 degrees of freedom
```

```
Multiple R-Squared: 0.9349, Adjusted R-squared: 0.9345
```

```
F-statistic: 2617 on 3 and 547 degrees of freedom, p-value: 0
```

```
> #actually it was unnecessary to construct the contrast matrix.
> #If we regress with a categorical variable without
> #an intercept, R will automatically use dummies for all levels:
> lm(lnwage ~ union + ed - 1, data=cps85)
```

```
Call:
```

```
lm(formula = lnwage ~ union + ed - 1, data = cps85)
```

```
Coefficients:
```

```
unionNonunion  unionUnion      ed
      0.9926      1.2909      0.0778
```

```
> ##Exercise 2d (1 point) Why is it not possible to include two dummies plus
> # an intercept? Because the two dummies add to 1,
> # you have perfect collinearity
```

```
> ###Exercise 3a (2 points):
```

```
> summary(lm(lnwage ~ ed + ex + I(ex^2), data=cps78))
```

```
> #All coefficients are highly significant, but the R^2 is only 0.2402
```

```
> #Returns to experience are positive and decline with increase in experience
```

```
> ##Exercise 3b (2 points):
```

```
> summary(lm(lnwage ~ gender + ed + ex + I(ex^2), data=cps78))
```

```
> contrasts(cps78$gender)
```

```
> #We see here that gender is coded 0 for female and 1 for male;
```

```
> #by default, the levels in a factor variable occur in alphabetical order
> #Intercept in our regression = 0.1909203 (this is for female),
> #genderMale has coefficient = 0.3351771,
> #i.e., the intercept for women is 0.5260974
> #Gender is highly significant
> ##Exercise 3c (2 points):
> summary(lm(lnwage ~ gender + marr + ed + ex + I(ex^2), data=cps78))
> #Coefficient of marr in this is insignificant
> ##Exercise 3d (1 point) asks to construct a variable which we do not
> #not need when we use factor variables
> ##Exercise 3e (3 points): For interaction term do
> summary(lm(lnwage ~ gender * marr + ed + ex + I(ex^2), data=cps78))
```

```
Call:
```

```
lm(formula = lnwage ~ gender * marr + ed + ex + I(ex^2), data = cps78)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-2.45524 -0.24566  0.01969  0.23102  1.42437
```

```
Coefficients:
```

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    0.1893919  0.1042613   1.817  0.06984 .
genderMale     0.3908782  0.0467018  8.370 4.44e-16 ***
marrSingle     0.0507811  0.0557198   0.911  0.36251
ed             0.0738640  0.0066154 11.165 < 2e-16 ***
ex             0.0265297  0.0049741  5.334 1.42e-07 ***
I(ex^2)       -0.0003161  0.0001057 -2.990  0.00291 **
genderMale:marrSingle -0.1586452  0.0750830 -2.113  0.03506 *
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.3959 on 543 degrees of freedom
```

```
Multiple R-Squared: 0.3547, Adjusted R-squared: 0.3476
```

```
F-statistic: 49.75 on 6 and 543 degrees of freedom, p-value: 0
```

```
> #Being married raises the wage for men by 13% but lowers it for women
```

```
> ##Exercise 4a (5 points):
```

```
> summary(lm(lnwage ~ union + gender + race + ed + ex + I(ex^2), data=cps78))
```

```
Call:
lm(formula = lnwage ~ union + gender + race + ed + ex + I(ex^2),
    data = cps78)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-2.41914 -0.23674  0.01682  0.21821  1.31584
```

```
Coefficients:
```

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.1549723  0.1068589   1.450  0.14757
unionUnion   0.2071429  0.0368503   5.621 3.04e-08 ***
genderMale   0.3060477  0.0344415   8.886 < 2e-16 ***
raceNonwh   -0.1301175  0.0830156  -1.567  0.11761
raceOther    0.0271477  0.0688277   0.394  0.69342
ed           0.0746097  0.0066521  11.216 < 2e-16 ***
ex           0.0261914  0.0047174   5.552 4.43e-08 ***
I(ex^2)     -0.0003082  0.0001015  -3.035  0.00252 **
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.3845 on 542 degrees of freedom
```

```
Multiple R-Squared: 0.3924, Adjusted R-squared: 0.3846
```

```
F-statistic: 50.01 on 7 and 542 degrees of freedom, p-value: 0
```

```
> exp(-0.1301175)
```

```
[1] 0.8779923
```

```
> #Being Hispanic lowers wages by 2.7%, byut being black lowers them
```

```
> #by 12.2 %
```

```
> #At what level of ex is lnwage maximized?
```

```
> #effect = 0.0261914 * ex -0.0003082 * ex^2
```

```
> #derivative = 0.0261914 - 2 * 0.0003082 * ex
```

```
> #derivative = 0 for ex=0.0261914/(2*0.0003082)
```

```
> 0.0261914/(2*0.0003082)
```

```
[1] 42.49091
```

```
> # age - ed - 6 = 42.49091
```

```
> # age = ed + 48.49091
> # for 8, 12, and 16 years of schooling the max earnings
> # are at ages 56.5, 60.5, and 64.5 years
> ##Exercise 4b (4 points) is a graph, not done here
> ##Exercise 4c (5 points)
> summary(lm(lnwage ~ gender + union + race + ed + ex + I(ex^2) + I
```

```
Call:
```

```
lm(formula = lnwage ~ gender + union + race + ed + ex + I(ex^2) +
    I(ed * ex), data = cps78)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-2.41207 -0.23922  0.01463  0.21645  1.32051
```

```
Coefficients:
```

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.0396495  0.1789073   0.222  0.824693
genderMale   0.3042639  0.0345241   8.813 < 2e-16 ***
unionUnion   0.2074045  0.0368638   5.626 2.96e-08 ***
raceNonwh   -0.1323898  0.0830908  -1.593  0.111673
raceOther    0.0319829  0.0691124   0.463  0.643718
ed           0.0824154  0.0117716   7.001 7.55e-12 ***
ex           0.0328854  0.0095716   3.436 0.000636 ***
I(ex^2)     -0.0003574  0.0001186  -3.013  0.002704 **
I(ed * ex)  -0.0003813  0.0004744  -0.804  0.421835
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.3846 on 541 degrees of freedom
```

```
Multiple R-Squared: 0.3932, Adjusted R-squared: 0.3842
```

```
F-statistic: 43.81 on 8 and 541 degrees of freedom, p-value: 0
```

```
> #Maximum earnings ages must be computed as before
```

```
> ##Exercise 4d (4 points) not done here
```

```
> ##Exercise 4e (6 points) not done here
```

```
> ###Exercise 5a (3 points):
```

```
> #Naive approach to estimate impact of unionization on wages:
```

```
> summary(lm(lnwage ~ gender + union + race + ed + ex + I(ex^2), dat
```

```
Call:
lm(formula = lnwage ~ gender + union + race + ed + ex + I(ex^2),
    data = cps78)
```

Residuals:

Min	1Q	Median	3Q	Max
-2.41914	-0.23674	0.01682	0.21821	1.31584

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.1549723	0.1068589	1.450	0.14757
genderMale	0.3060477	0.0344415	8.886	< 2e-16 ***
unionUnion	0.2071429	0.0368503	5.621	3.04e-08 ***
raceNonwh	-0.1301175	0.0830156	-1.567	0.11761
raceOther	0.0271477	0.0688277	0.394	0.69342
ed	0.0746097	0.0066521	11.216	< 2e-16 ***
ex	0.0261914	0.0047174	5.552	4.43e-08 ***
I(ex^2)	-0.0003082	0.0001015	-3.035	0.00252 **

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3845 on 542 degrees of freedom  
Multiple R-Squared: 0.3924, Adjusted R-squared: 0.3846  
F-statistic: 50.01 on 7 and 542 degrees of freedom, p-value: 0

```
> # What is wrong with the above? It assumes that unions
> # only affect the intercept, everything else is the same
> ##Exercise 5b (2 points)
```

```
> tapply(lnwage, union, mean)
```

Nonun	Union
1.600901	1.863137

```
> tapply(ed, union, mean)
```

Nonun	Union
12.76178	12.02381

```
> table(gender, union)
```

	union	
gender	Nonun	Union
Female	159	48

```
Male      223   120
```

```
> table(race, union)
```

	union	
race	Nonun	Union
Hisp	29	7
Nonwh	35	22
Other	318	139

```
> 7/(7+29)
```

```
[1] 0.1944444
```

```
> 22/(22+35)
```

```
[1] 0.3859649
```

```
> 139/(318+139)
```

```
[1] 0.3041575
```

```
> #19% of Hispanic, 39% of Nonwhite, and 30% of other (white) workers
```

```
> #in the sample are in unions
```

```
> ##Exercise 5c (3 points)
```

```
> summary(lm(lnwage ~ gender + race + ed + ex + I(ex^2), data=cps78))
```

Call:

```
lm(formula = lnwage ~ gender + race + ed + ex + I(ex^2), data = cps78,
    subset = union == "Union")
```

Residuals:

Min	1Q	Median	3Q	Max
-2.3307	-0.1853	0.0160	0.2199	1.1992

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.9261456	0.2321964	3.989	0.000101 ***
genderMale	0.2239370	0.0684894	3.270	0.001317 **
raceNonwh	-0.3066717	0.1742287	-1.760	0.080278 .
raceOther	-0.0741660	0.1562131	-0.475	0.635591
ed	0.0399500	0.0138311	2.888	0.004405 **
ex	0.0313820	0.0098938	3.172	0.001814 **
I(ex^2)	-0.0004526	0.0002022	-2.239	0.026535 *

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3928 on 161 degrees of freedom

Multiple R-Squared: 0.2019, Adjusted R-squared: 0.1721

F-statistic: 6.787 on 6 and 161 degrees of freedom, p-value: 1.975e-06

```
> summary(lm(lnwage ~ gender + race + ed + ex + I(ex^2), data=cps78, subset = union == "Nonun"))
```

Call:

```
lm(formula = lnwage ~ gender + race + ed + ex + I(ex^2), data = cps78, subset = union == "Nonun")
```

Residuals:

Min	1Q	Median	3Q	Max
-1.39107	-0.23775	0.01040	0.23337	1.29073

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.0095668	0.1193399	-0.080	0.9361
genderMale	0.3257661	0.0397961	8.186	4.22e-15 ***
raceNonwh	-0.0652018	0.0960570	-0.679	0.4977
raceOther	0.0444133	0.0761628	0.583	0.5602
ed	0.0852212	0.0075554	11.279	< 2e-16 ***
ex	0.0253813	0.0053710	4.726	3.25e-06 ***
I(ex^2)	-0.0002841	0.0001187	-2.392	0.0172 *

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3778 on 375 degrees of freedom

Multiple R-Squared: 0.4229, Adjusted R-squared: 0.4137

F-statistic: 45.8 on 6 and 375 degrees of freedom, p-value: 0

```
> #Are union-nonunion differences larger for females than males?
```

```
> #For this look at the intercepts for males and females in
```

```
> #the two regressions. Say for white males and females:
```

```
> 0.9261456-0.0741660+0.2239370
```

```
[1] 1.075917
```

```
> 0.9261456-0.0741660
```

```
[1] 0.8519796
```

```
> -0.0095668+0.0444133+0.3257661
```

```
[1] 0.3606126
```

```
> -0.0095668+0.0444133
```

```
[1] 0.0348465
```

```
> 1.075917-0.3606126
```

```
[1] 0.7153044
```

```
> 0.8519796-0.0348465
```

```
[1] 0.8171331
```

```
>
```

```
> #White Males White Females
```

```
> #Union 1.075917 0.8519796
```

```
> #Nonunion 0.3606126 0.0348465
```

```
> #Difference 0.7153044 0.8171331
```

```
> #Difference is greater for women
```

```
> ###Exercise 6a (5 points)
```

```
> summary(lm(lnwage ~ gender + union + race + ed + ex + I(ex^2)))
```

Call:

```
lm(formula = lnwage ~ gender + union + race + ed + ex + I(ex^2))
```

Residuals:

Min	1Q	Median	3Q	Max
-2.41914	-0.23674	0.01682	0.21821	1.31584

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.1549723	0.1068589	1.450	0.14757
genderMale	0.3060477	0.0344415	8.886	< 2e-16 ***
unionUnion	0.2071429	0.0368503	5.621	3.04e-08 ***
raceNonwh	-0.1301175	0.0830156	-1.567	0.11761
raceOther	0.0271477	0.0688277	0.394	0.69342
ed	0.0746097	0.0066521	11.216	< 2e-16 ***
ex	0.0261914	0.0047174	5.552	4.43e-08 ***
I(ex^2)	-0.0003082	0.0001015	-3.035	0.00252 **

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3845 on 542 degrees of freedom

Multiple R-Squared: 0.3924, Adjusted R-squared: 0.3846

F-statistic: 50.01 on 7 and 542 degrees of freedom, p-value: 0

```
> #To test whether Nonwh and Hisp have same intercept
```

```
> #one might generate a contrast matrix which collapses those
```

```
> #two and then run it and make an F-test
```

```
> #or make a contrast matrix which has this difference as one of
> #the dummies and use the t-test for that dummy
> ##Exercise 6b (2 points)
> table(race)
race
  Hisp Nonwh Other
    36   57  457
> tapply(lnwage, race, mean)
  Hisp   Nonwh   Other
1.529647 1.513404 1.713829
> tapply(lnwage, race, ed)
Error in get(x, envir, mode, inherits) : variable "ed" was not found
> tapply(ed, race, mean)
  Hisp   Nonwh   Other
10.30556 11.71930 12.81400
> table(gender, race)
      race
gender  Hisp Nonwh Other
  Female   12   28  167
  Male    24   29  290
> #Blacks, almost as many women than men, hispanic twice as many men,
> #Whites in between

>
> #Additional stuff:
> #There are two outliers in cps78 with wages of less than $1 per hour,
> #Both service workers, perhaps waitresses who did not report her tips?
> #What are the commands for extracting certain observations
> #by certain criteria and just print them? The split command.
>
> #Interesting to do
> loess(lnwage ~ ed + ex, data=cps78)
> #loess is appropriate here because there are strong interation terms
> #How can one do loess after taking out the effects of gender for instan
> #Try the following, but I did not try it out yet:
> gam(lnwage ~ lo(ed,ex) + gender, data=cps78)
> #I should put more plotting commands in!
```

analog would be to find two vectors in a plane  $\hat{\phi}$  and  $\tilde{\phi}$ . In each component (i.e., projection on the axes),  $\hat{\phi}$  is closer to the origin than  $\tilde{\phi}$ . But in the projection on the diagonal,  $\tilde{\phi}$  is closer to the origin than  $\hat{\phi}$ .

ANSWER. In the simplest counterexample, all variables involved are constants:  $\phi = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ,  $\hat{\phi} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ , and  $\tilde{\phi} = \begin{bmatrix} -2 \\ 2 \end{bmatrix}$ .

## CHAPTER 17

# The Mean Squared Error as an Initial Criterion of Precision

The question how “close” two random variables are to each other is a central concern in statistics. The goal of statistics is to find observed random variables which are “close” to the unobserved parameters or random outcomes of interest. These observed random variables are usually called “estimators” if the unobserved magnitude is nonrandom, and “predictors” if it is random. For *scalar* random variables we will use the mean squared error as a criterion for closeness. Its definition is  $\text{MSE}[\hat{\phi}; \phi]$  (read it: mean squared error of  $\hat{\phi}$  as an estimator or predictor, whatever the case may be, of  $\phi$ ):

$$(17.0.1) \quad \text{MSE}[\hat{\phi}; \phi] = \text{E}[(\hat{\phi} - \phi)^2]$$

For our purposes, therefore, the estimator (or predictor)  $\hat{\phi}$  of the unknown parameter (or unobserved random variable)  $\phi$  is no worse than the alternative  $\tilde{\phi}$  if  $\text{MSE}[\hat{\phi}; \phi] \leq \text{MSE}[\tilde{\phi}; \phi]$ . This is a criterion which can be applied before any observations are collected and actual estimations are made; it is an “initial” criterion regarding the expected average performance in a series of future trials (even though, in economics, usually only one trial is made).

### 17.1. Comparison of Two Vector Estimators

If one wants to compare two *vector* estimators, say  $\hat{\phi}$  and  $\tilde{\phi}$ , it is often impossible to say which of two estimators is better. It may be the case that  $\hat{\phi}_1$  is better than  $\tilde{\phi}_1$  (in terms of MSE or some other criterion), but  $\hat{\phi}_2$  is worse than  $\tilde{\phi}_2$ . And even if every component  $\phi_i$  is estimated better by  $\hat{\phi}_i$  than by  $\tilde{\phi}_i$ , certain linear combinations  $\mathbf{t}^\top \phi$  of the components of  $\phi$  may be estimated better by  $\mathbf{t}^\top \tilde{\phi}$  than by  $\mathbf{t}^\top \hat{\phi}$ .

PROBLEM 240. *2 points* Construct an example of two vector estimators  $\hat{\phi}$  and  $\tilde{\phi}$  of the same random vector  $\phi = [\phi_1 \ \phi_2]^\top$ , so that  $\text{MSE}[\hat{\phi}_i; \phi_i] < \text{MSE}[\tilde{\phi}_i; \phi_i]$  for  $i = 1, 2$  but  $\text{MSE}[\hat{\phi}_1 + \hat{\phi}_2; \phi_1 + \phi_2] > \text{MSE}[\tilde{\phi}_1 + \tilde{\phi}_2; \phi_1 + \phi_2]$ . *Hint: it is easiest to use an example in which all random variables are constants. Another hint: the geometric*

One can only then say unambiguously that the vector  $\hat{\phi}$  is a no worse estimator than  $\tilde{\phi}$  if its MSE is smaller or equal for every linear combination. Theorem 17.1.1 will show that this is the case if and only if the *MSE-matrix* of  $\hat{\phi}$  is smaller, by nonnegative definite matrix, than that of  $\tilde{\phi}$ . If this is so, then theorem 17.1.1 says that not only the MSE of all linear transformations, but also all other nonnegative definite quadratic loss functions involving these vectors (such as the trace of the *MSE-matrix*, which is an often-used criterion) are minimized. In order to formulate and prove this, we first need a formal definition of the *MSE-matrix*. We write  $\mathcal{MSE}[\hat{\phi}; \phi]$  for the matrix and MSE for the scalar mean squared error. The *MSE-matrix* of  $\hat{\phi}$  as an estimator of  $\phi$  is defined as

$$(17.1.1) \quad \mathcal{MSE}[\hat{\phi}; \phi] = \mathcal{E}[(\hat{\phi} - \phi)(\hat{\phi} - \phi)^\top].$$

PROBLEM 241. *2 points* Let  $\theta$  be a vector of possibly random parameters, and  $\hat{\theta}$  an estimator of  $\theta$ . Show that

$$(17.1.2) \quad \mathcal{MSE}[\hat{\theta}; \theta] = \mathcal{V}[\hat{\theta} - \theta] + (\mathcal{E}[\hat{\theta} - \theta])(\mathcal{E}[\hat{\theta} - \theta])^\top.$$

*Don't assume the scalar result but make a proof that is good for vectors and scalars.*

ANSWER. For any random vector  $\mathbf{x}$  follows

$$\begin{aligned} \mathcal{E}[\mathbf{x}\mathbf{x}^\top] &= \mathcal{E}[(\mathbf{x} - \mathcal{E}[\mathbf{x}] + \mathcal{E}[\mathbf{x}])(\mathbf{x} - \mathcal{E}[\mathbf{x}] + \mathcal{E}[\mathbf{x}])^\top] \\ &= \mathcal{E}[(\mathbf{x} - \mathcal{E}[\mathbf{x}])(\mathbf{x} - \mathcal{E}[\mathbf{x}])^\top] - \mathcal{E}[(\mathbf{x} - \mathcal{E}[\mathbf{x}])\mathcal{E}[\mathbf{x}]^\top] - \mathcal{E}[\mathcal{E}[\mathbf{x}](\mathbf{x} - \mathcal{E}[\mathbf{x}])^\top] + \mathcal{E}[\mathcal{E}[\mathbf{x}]\mathcal{E}[\mathbf{x}]^\top] \\ &= \mathcal{V}[\mathbf{x}] - \mathbf{O} - \mathbf{O} + \mathcal{E}[\mathbf{x}]\mathcal{E}[\mathbf{x}]^\top. \end{aligned}$$

Setting  $\mathbf{x} = \hat{\theta} - \theta$  the statement follows.

If  $\theta$  is nonrandom, formula (17.1.2) simplifies slightly, since in this case  $\mathcal{V}[\hat{\theta} - \theta] = \mathcal{V}[\hat{\theta}]$ . In this case, the *MSE matrix* is the covariance matrix plus the squared bias matrix. If  $\theta$  is nonrandom and in addition  $\hat{\theta}$  is unbiased, then the *MSE-matrix* coincides with the covariance matrix.

THEOREM 17.1.1. *Assume  $\hat{\phi}$  and  $\tilde{\phi}$  are two estimators of the parameter  $\phi$  (which is allowed to be random itself). Then conditions (17.1.3), (17.1.4), and*

(17.1.5) are equivalent:

$$(17.1.3) \quad \text{For every constant vector } \mathbf{t}, \quad \text{MSE}[\mathbf{t}^\top \hat{\boldsymbol{\phi}}; \mathbf{t}^\top \boldsymbol{\phi}] \leq \text{MSE}[\mathbf{t}^\top \tilde{\boldsymbol{\phi}}; \mathbf{t}^\top \boldsymbol{\phi}]$$

$$(17.1.4) \quad \mathcal{MSE}[\tilde{\boldsymbol{\phi}}; \boldsymbol{\phi}] - \mathcal{MSE}[\hat{\boldsymbol{\phi}}; \boldsymbol{\phi}] \quad \text{is a nonnegative definite matrix}$$

$$(17.1.5) \quad \text{For every nnd } \boldsymbol{\Theta}, \quad \text{E}[(\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top \boldsymbol{\Theta} (\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})] \leq \text{E}[(\tilde{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top \boldsymbol{\Theta} (\tilde{\boldsymbol{\phi}} - \boldsymbol{\phi})].$$

PROOF. Call  $\mathcal{MSE}[\tilde{\boldsymbol{\phi}}; \boldsymbol{\phi}] = \sigma^2 \boldsymbol{\Xi}$  and  $\mathcal{MSE}[\hat{\boldsymbol{\phi}}; \boldsymbol{\phi}] = \sigma^2 \boldsymbol{\Omega}$ . To show that (17.1.3) implies (17.1.4), simply note that  $\text{MSE}[\mathbf{t}^\top \hat{\boldsymbol{\phi}}; \mathbf{t}^\top \boldsymbol{\phi}] = \sigma^2 \mathbf{t}^\top \boldsymbol{\Omega} \mathbf{t}$  and likewise  $\text{MSE}[\mathbf{t}^\top \tilde{\boldsymbol{\phi}}; \mathbf{t}^\top \boldsymbol{\phi}] = \sigma^2 \mathbf{t}^\top \boldsymbol{\Xi} \mathbf{t}$ . Therefore (17.1.3) is equivalent to  $\mathbf{t}^\top (\boldsymbol{\Xi} - \boldsymbol{\Omega}) \mathbf{t} \geq 0$  for all  $\mathbf{t}$ , which is the defining property making  $\boldsymbol{\Xi} - \boldsymbol{\Omega}$  nonnegative definite.

Here is the proof that (17.1.4) implies (17.1.5):

$$\begin{aligned} \text{E}[(\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top \boldsymbol{\Theta} (\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})] &= \text{E}[\text{tr}((\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top \boldsymbol{\Theta} (\hat{\boldsymbol{\phi}} - \boldsymbol{\phi}))] = \\ &= \text{E}[\text{tr}(\boldsymbol{\Theta} (\hat{\boldsymbol{\phi}} - \boldsymbol{\phi}) (\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top)] = \text{tr}(\boldsymbol{\Theta} \text{E}[(\hat{\boldsymbol{\phi}} - \boldsymbol{\phi}) (\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top]) = \sigma^2 \text{tr}(\boldsymbol{\Theta} \boldsymbol{\Omega}) \end{aligned}$$

and in the same way

$$\text{E}[(\tilde{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top \boldsymbol{\Theta} (\tilde{\boldsymbol{\phi}} - \boldsymbol{\phi})] = \sigma^2 \text{tr}(\boldsymbol{\Theta} \boldsymbol{\Xi}).$$

The difference in the expected quadratic forms is therefore  $\sigma^2 \text{tr}(\boldsymbol{\Theta} (\boldsymbol{\Xi} - \boldsymbol{\Omega}))$ . By assumption,  $\boldsymbol{\Xi} - \boldsymbol{\Omega}$  is nonnegative definite. Therefore, by theorem A.5.6 in the Mathematical Appendix, or by Problem 242 below, this trace is nonnegative.

To complete the proof, (17.1.5) has (17.1.3) as a special case if one sets  $\boldsymbol{\Theta} = \mathbf{t} \mathbf{t}^\top$ .  $\square$

PROBLEM 242. Show that if  $\boldsymbol{\Theta}$  and  $\boldsymbol{\Sigma}$  are symmetric and nonnegative definite, then  $\text{tr}(\boldsymbol{\Theta} \boldsymbol{\Sigma}) \geq 0$ . You are allowed to use that  $\text{tr}(\mathbf{A} \mathbf{B}) = \text{tr}(\mathbf{B} \mathbf{A})$ , that the trace of a nonnegative definite matrix is  $\geq 0$ , and Problem 118 (which is trivial).

$$\text{ANSWER. Write } \boldsymbol{\Theta} = \mathbf{R} \mathbf{R}^\top; \text{ then } \text{tr}(\boldsymbol{\Theta} \boldsymbol{\Sigma}) = \text{tr}(\mathbf{R} \mathbf{R}^\top \boldsymbol{\Sigma}) = \text{tr}(\mathbf{R}^\top \boldsymbol{\Sigma} \mathbf{R}) \geq 0. \quad \square$$

PROBLEM 243. Consider two very simple-minded estimators of the unknown nonrandom parameter vector  $\boldsymbol{\phi} = \begin{bmatrix} \phi_1 \\ \phi_2 \end{bmatrix}$ . Neither of these estimators depends on any observations, they are constants. The first estimator is  $\hat{\boldsymbol{\phi}} = \begin{bmatrix} 11 \\ 11 \end{bmatrix}$ , and the second is  $\tilde{\boldsymbol{\phi}} = \begin{bmatrix} 12 \\ 8 \end{bmatrix}$ .

• a. 2 points Compute the  $\mathcal{MSE}$ -matrices of these two estimators if the true value of the parameter vector is  $\boldsymbol{\phi} = \begin{bmatrix} 10 \\ 10 \end{bmatrix}$ . For which estimator is the trace of the  $\mathcal{MSE}$  matrix smaller?

ANSWER.  $\hat{\boldsymbol{\phi}}$  has smaller trace of the  $\mathcal{MSE}$ -matrix.

$$\begin{aligned} \hat{\boldsymbol{\phi}} - \boldsymbol{\phi} &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ \mathcal{MSE}[\hat{\boldsymbol{\phi}}; \boldsymbol{\phi}] &= \text{E}[(\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})(\hat{\boldsymbol{\phi}} - \boldsymbol{\phi})^\top] \\ &= \text{E}\left[\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix}\right] = \text{E}\left[\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}\right] = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \\ \tilde{\boldsymbol{\phi}} - \boldsymbol{\phi} &= \begin{bmatrix} 2 \\ -2 \end{bmatrix} \\ \mathcal{MSE}[\tilde{\boldsymbol{\phi}}; \boldsymbol{\phi}] &= \begin{bmatrix} 4 & -4 \\ -4 & 4 \end{bmatrix} \end{aligned}$$

Note that both  $\mathcal{MSE}$ -matrices are singular, i.e., both estimators allow an error-free look at certain linear combinations of the parameter vector.

• b. 1 point Give two vectors  $\mathbf{g} = \begin{bmatrix} g_1 \\ g_2 \end{bmatrix}$  and  $\mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}$  satisfying  $\text{MSE}[\mathbf{g}^\top \hat{\boldsymbol{\phi}}; \mathbf{g}^\top \boldsymbol{\phi}] < \text{MSE}[\mathbf{g}^\top \tilde{\boldsymbol{\phi}}; \mathbf{g}^\top \boldsymbol{\phi}]$  and  $\text{MSE}[\mathbf{h}^\top \hat{\boldsymbol{\phi}}; \mathbf{h}^\top \boldsymbol{\phi}] > \text{MSE}[\mathbf{h}^\top \tilde{\boldsymbol{\phi}}; \mathbf{h}^\top \boldsymbol{\phi}]$  ( $\mathbf{g}$  and  $\mathbf{h}$  are not unique, there are many possibilities).

ANSWER. With  $\mathbf{g} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$  and  $\mathbf{h} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  for instance we get  $\mathbf{g}^\top \hat{\boldsymbol{\phi}} - \mathbf{g}^\top \boldsymbol{\phi} = 0$ ,  $\mathbf{g}^\top \tilde{\boldsymbol{\phi}} - \mathbf{g}^\top \boldsymbol{\phi} = 0$ ,  $\mathbf{h}^\top \hat{\boldsymbol{\phi}} - \mathbf{h}^\top \boldsymbol{\phi} = 2$ ,  $\mathbf{h}^\top \tilde{\boldsymbol{\phi}} - \mathbf{h}^\top \boldsymbol{\phi} = 0$ , therefore  $\text{MSE}[\mathbf{g}^\top \hat{\boldsymbol{\phi}}; \mathbf{g}^\top \boldsymbol{\phi}] = 0$ ,  $\text{MSE}[\mathbf{g}^\top \tilde{\boldsymbol{\phi}}; \mathbf{g}^\top \boldsymbol{\phi}] = 0$ ,  $\text{MSE}[\mathbf{h}^\top \hat{\boldsymbol{\phi}}; \mathbf{h}^\top \boldsymbol{\phi}] = 4$ ,  $\text{MSE}[\mathbf{h}^\top \tilde{\boldsymbol{\phi}}; \mathbf{h}^\top \boldsymbol{\phi}] = 0$ . An alternative way to compute this is e.g.

$$\mathcal{MSE}[\mathbf{h}^\top \tilde{\boldsymbol{\phi}}; \mathbf{h}^\top \boldsymbol{\phi}] = \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} 4 & -4 \\ -4 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 16$$

• c. 1 point Show that neither  $\mathcal{MSE}[\hat{\boldsymbol{\phi}}; \boldsymbol{\phi}] - \mathcal{MSE}[\tilde{\boldsymbol{\phi}}; \boldsymbol{\phi}]$  nor  $\mathcal{MSE}[\tilde{\boldsymbol{\phi}}; \boldsymbol{\phi}] - \mathcal{MSE}[\hat{\boldsymbol{\phi}}; \boldsymbol{\phi}]$  is a nonnegative definite matrix. Hint: you are allowed to use the mathematical fact that if a matrix is nonnegative definite, then its determinant is nonnegative.

ANSWER.

$$(17.1.6) \quad \mathcal{MSE}[\tilde{\boldsymbol{\phi}}; \boldsymbol{\phi}] - \mathcal{MSE}[\hat{\boldsymbol{\phi}}; \boldsymbol{\phi}] = \begin{bmatrix} 3 & -5 \\ -5 & 3 \end{bmatrix}$$

Its determinant is negative, and the determinant of its negative is also negative.

## CHAPTER 18

## Sampling Properties of the Least Squares Estimator

The estimator  $\hat{\beta}$  was derived from a *geometric* argument, and everything which we showed so far are what [DM93, p. 3] calls its *numerical* as opposed to its *statistical* properties. But  $\hat{\beta}$  has also nice *statistical* or *sampling* properties. We are assuming right now the specification given in (14.1.3), in which  $\mathbf{X}$  is an arbitrary matrix of full column rank, and we are not assuming that the errors must be Normally distributed. The assumption that  $\mathbf{X}$  is nonrandom means that repeated samples are taken with the same  $\mathbf{X}$ -matrix. This is often true for experimental data, but not in econometrics. The sampling properties which we are really interested in are those where also the  $\mathbf{X}$ -matrix is random; we will derive those later. For this later derivation, the properties with fixed  $\mathbf{X}$ -matrix, which we are going to discuss presently, will be needed as an intermediate step. The assumption of fixed  $\mathbf{X}$  is therefore a preliminary technical assumption, to be dropped later.

In order to know how good the estimator  $\hat{\beta}$  is, one needs the statistical properties of its “sampling error”  $\hat{\beta} - \beta$ . This sampling error has the following formula:

$$\begin{aligned} \hat{\beta} - \beta &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \beta = \\ (18.0.7) \quad &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \beta) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} \end{aligned}$$

From (18.0.7) follows immediately that  $\hat{\beta}$  is unbiased, since  $\mathcal{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}] = \mathbf{o}$ .

Unbiasedness does not make an estimator better, but many good estimators are unbiased, and it simplifies the math.

We will use the  $\mathcal{MSE}$ -matrix as a criterion for how good an estimator of a vector of unobserved parameters is. Chapter 17 gave some reasons why this is a sensible criterion (compare [DM93, Chapter 5.5]).

### 18.1. The Gauss Markov Theorem

Returning to the least squares estimator  $\hat{\beta}$ , one obtains, using (18.0.7), that

$$\begin{aligned} \mathcal{MSE}[\hat{\beta}; \beta] &= \mathcal{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top] = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathcal{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top] \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \\ (18.1.1) \quad &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned}$$

This is a very simple formula. Its most interesting aspect is that this  $\mathcal{MSE}$  matrix does not depend on the value of the true  $\beta$ . In particular this means that it is *bounded* with respect to  $\beta$ , which is important for someone who wants to be assured of a certain accuracy even in the worst possible situation.

**PROBLEM 244.** 2 points Compute the  $\mathcal{MSE}$ -matrix  $\mathcal{MSE}[\hat{\boldsymbol{\varepsilon}}; \boldsymbol{\varepsilon}] = \mathcal{E}[(\hat{\boldsymbol{\varepsilon}} - \boldsymbol{\varepsilon})(\hat{\boldsymbol{\varepsilon}} - \boldsymbol{\varepsilon})^\top]$  of the residuals as predictors of the disturbances.

**ANSWER.** Write  $\hat{\boldsymbol{\varepsilon}} - \boldsymbol{\varepsilon} = \mathbf{M} \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon} = (\mathbf{M} - \mathbf{I}) \boldsymbol{\varepsilon} = -\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}$ ; therefore  $\mathcal{MSE}[\hat{\boldsymbol{\varepsilon}}; \boldsymbol{\varepsilon}] = \mathcal{E}[\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}] = \sigma^2 \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ . Alternatively, start with  $\hat{\boldsymbol{\varepsilon}} - \boldsymbol{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}} - \boldsymbol{\varepsilon} = \mathbf{X} \beta - \hat{\mathbf{y}} - \boldsymbol{\varepsilon} = \mathbf{X} (\beta - \hat{\beta})$ . This allows to use  $\mathcal{MSE}[\hat{\boldsymbol{\varepsilon}}; \boldsymbol{\varepsilon}] = \mathbf{X} \mathcal{MSE}[\hat{\beta}; \beta] \mathbf{X}^\top = \sigma^2 \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .

**PROBLEM 245.** 2 points Let  $\mathbf{v}$  be a random vector that is a linear transformation of  $\mathbf{y}$ , i.e.,  $\mathbf{v} = \mathbf{T} \mathbf{y}$  for some constant matrix  $\mathbf{T}$ . Furthermore  $\mathbf{v}$  satisfies  $\mathcal{E}[\mathbf{v}] = \mathbf{o}$ . Show that from this follows  $\mathbf{v} = \mathbf{T} \hat{\boldsymbol{\varepsilon}}$ . (In other words, no other transformation of  $\mathbf{y}$  with zero expected value is more “comprehensive” than  $\boldsymbol{\varepsilon}$ . However there are many other transformations of  $\mathbf{y}$  with zero expected value which are as “comprehensive” as  $\boldsymbol{\varepsilon}$ .)

**ANSWER.**  $\mathcal{E}[\mathbf{v}] = \mathbf{T} \mathbf{X} \beta$  must be  $\mathbf{o}$  whatever the value of  $\beta$ . Therefore  $\mathbf{T} \mathbf{X} = \mathbf{O}$ , from which follows  $\mathbf{T} \mathbf{M} = \mathbf{T}$ . Since  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M} \mathbf{y}$ , this gives immediately  $\mathbf{v} = \mathbf{T} \hat{\boldsymbol{\varepsilon}}$ . (This is the statistical implication of the mathematical fact that  $\mathbf{M}$  is a deficiency matrix of  $\mathbf{X}$ .)

**PROBLEM 246.** 2 points Show that  $\hat{\beta}$  and  $\hat{\boldsymbol{\varepsilon}}$  are uncorrelated, i.e.,  $\text{cov}[\hat{\beta}_i, \hat{\boldsymbol{\varepsilon}}_j] = 0$  for all  $i, j$ . Defining the covariance matrix  $\mathcal{C}[\hat{\beta}, \hat{\boldsymbol{\varepsilon}}]$  as that matrix whose  $(i, j)$ -element is  $\text{cov}[\hat{\beta}_i, \hat{\boldsymbol{\varepsilon}}_j]$ , this can also be written as  $\mathcal{C}[\hat{\beta}, \hat{\boldsymbol{\varepsilon}}] = \mathbf{O}$ . Hint: The covariance matrix satisfies the rules  $\mathcal{C}[\mathbf{A} \mathbf{y}, \mathbf{B} \mathbf{z}] = \mathbf{A} \mathcal{C}[\mathbf{y}, \mathbf{z}] \mathbf{B}^\top$  and  $\mathcal{C}[\mathbf{y}, \mathbf{y}] = \mathcal{V}[\mathbf{y}]$ . (Other rules for the covariance matrix, which will not be needed here, are  $\mathcal{C}[\mathbf{z}, \mathbf{y}] = (\mathcal{C}[\mathbf{y}, \mathbf{z}])^\top$ ,  $\mathcal{C}[\mathbf{x} + \mathbf{y}, \mathbf{z}] = \mathcal{C}[\mathbf{x}, \mathbf{z}] + \mathcal{C}[\mathbf{y}, \mathbf{z}]$ ,  $\mathcal{C}[\mathbf{x}, \mathbf{y} + \mathbf{z}] = \mathcal{C}[\mathbf{x}, \mathbf{y}] + \mathcal{C}[\mathbf{x}, \mathbf{z}]$ , and  $\mathcal{C}[\mathbf{y}, \mathbf{c}] = \mathbf{O}$  if  $\mathbf{c}$  is a vector of constants.)

**ANSWER.**  $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  and  $\mathbf{B} = \mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ , therefore  $\mathcal{C}[\hat{\beta}, \hat{\boldsymbol{\varepsilon}}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \mathbf{O}$ .

**PROBLEM 247.** 4 points Let  $\mathbf{y} = \mathbf{X} \beta + \boldsymbol{\varepsilon}$  be a regression model with intercept, which the first column of  $\mathbf{X}$  is the vector  $\mathbf{1}$ , and let  $\hat{\beta}$  the least squares estimator



$\beta$ . Show that the covariance matrix between  $\bar{\mathbf{y}}$  and  $\hat{\beta}$ , which is defined as the matrix (here consisting of one row only) that contains all the covariances

$$(18.1.2) \quad \mathcal{C}[\bar{\mathbf{y}}, \hat{\beta}] \equiv [\text{cov}[\bar{\mathbf{y}}, \hat{\beta}_1] \quad \text{cov}[\bar{\mathbf{y}}, \hat{\beta}_2] \quad \cdots \quad \text{cov}[\bar{\mathbf{y}}, \hat{\beta}_k]]$$

has the following form:  $\mathcal{C}[\bar{\mathbf{y}}, \hat{\beta}] = \frac{\sigma^2}{n} [1 \quad 0 \quad \cdots \quad 0]$  where  $n$  is the number of observations. Hint: That the regression has an intercept term as first column of the  $\mathbf{X}$ -matrix means that  $\mathbf{X}\mathbf{e}^{(1)} = \mathbf{1}$ , where  $\mathbf{e}^{(1)}$  is the unit vector having 1 in the first place and zeros elsewhere, and  $\mathbf{1}$  is the vector which has ones everywhere.

ANSWER. Write both  $\bar{\mathbf{y}}$  and  $\hat{\beta}$  in terms of  $\mathbf{y}$ , i.e.,  $\bar{\mathbf{y}} = \frac{1}{n}\mathbf{1}^\top \mathbf{y}$  and  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ . Therefore

$$(18.1.3) \quad \mathcal{C}[\bar{\mathbf{y}}, \hat{\beta}] = \frac{1}{n} \mathbf{1}^\top \nu[\mathbf{y}] \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \frac{\sigma^2}{n} \mathbf{1}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \frac{\sigma^2}{n} \mathbf{e}^{(1)\top} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \frac{\sigma^2}{n} \mathbf{e}^{(1)\top}. \quad \square$$

**THEOREM 18.1.1.** Gauss-Markov Theorem:  $\hat{\beta}$  is the BLUE (Best Linear Unbiased Estimator) of  $\beta$  in the following vector sense: for every nonrandom coefficient vector  $\mathbf{t}$ ,  $\mathbf{t}^\top \hat{\beta}$  is the scalar BLUE of  $\mathbf{t}^\top \beta$ , i.e., every other linear unbiased estimator  $\tilde{\phi} = \mathbf{a}^\top \mathbf{y}$  of  $\phi = \mathbf{t}^\top \beta$  has a bigger MSE than  $\mathbf{t}^\top \hat{\beta}$ .

PROOF. Write the alternative linear estimator  $\tilde{\phi} = \mathbf{a}^\top \mathbf{y}$  in the form

$$(18.1.4) \quad \tilde{\phi} = (\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) \mathbf{y}$$

then the sampling error is

$$(18.1.5) \quad \begin{aligned} \tilde{\phi} - \phi &= (\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) (\mathbf{X}\beta + \boldsymbol{\varepsilon}) - \mathbf{t}^\top \beta \\ &= (\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) \boldsymbol{\varepsilon} + \mathbf{c}^\top \mathbf{X}\beta. \end{aligned}$$

By assumption, the alternative estimator is unbiased, i.e., the expected value of this sampling error is zero regardless of the value of  $\beta$ . This is only possible if  $\mathbf{c}^\top \mathbf{X} = \mathbf{o}^\top$ . But then it follows

$$\begin{aligned} \text{MSE}[\tilde{\phi}; \phi] &= \text{E}[(\tilde{\phi} - \phi)^2] = \text{E}[(\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top (\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t} + \mathbf{c})] = \\ &= \sigma^2 (\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) (\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t} + \mathbf{c}) = \sigma^2 \mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t} + \sigma^2 \mathbf{c}^\top \mathbf{c}, \end{aligned}$$

Here we needed again  $\mathbf{c}^\top \mathbf{X} = \mathbf{o}^\top$ . Clearly, this is minimized if  $\mathbf{c} = \mathbf{o}$ , in which case  $\tilde{\phi} = \mathbf{t}^\top \hat{\beta}$ .  $\square$

**PROBLEM 248.** 4 points Show: If  $\tilde{\beta}$  is a linear unbiased estimator of  $\beta$  and  $\hat{\beta}$  is the OLS estimator, then the difference of the MSE-matrices  $\text{MSE}[\tilde{\beta}; \beta] - \text{MSE}[\hat{\beta}; \beta]$  is nonnegative definite.

ANSWER. (Compare [DM93, p. 159].) Any other linear estimator  $\tilde{\beta}$  of  $\beta$  can be written as  $\tilde{\beta} = ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{C}) \mathbf{y}$ . Its expected value is  $\mathcal{E}[\tilde{\beta}] = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\beta + \mathbf{C}\mathbf{X}\beta$ .  $\tilde{\beta}$  to be unbiased, regardless of the value of  $\beta$ ,  $\mathbf{C}$  must satisfy  $\mathbf{C}\mathbf{X} = \mathbf{O}$ . But then it follows  $\text{MSE}[\tilde{\beta}; \beta] = \nu[\tilde{\beta}] = \sigma^2 ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{C}) (\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{C}^\top) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} + \sigma^2 \mathbf{C}\mathbf{C}^\top$  it exceeds the MSE-matrix of  $\hat{\beta}$  by a nonnegative definite matrix.

## 18.2. Digression about Minimax Estimators

Theorem 18.1.1 is a somewhat puzzling property of the least squares estimator since there is no reason in the world to restrict one's search for good estimators to unbiased estimators. An alternative and more enlightening characterization of  $\hat{\beta}$  does not use the concept of unbiasedness but that of a minimax estimator with respect to the MSE. For this I am proposing the following definition:

**DEFINITION 18.2.1.**  $\hat{\phi}$  is the linear minimax estimator of the scalar parameter  $\phi$  with respect to the MSE if and only if for every other linear estimator  $\tilde{\phi}$  there exists a value of the parameter vector  $\beta_0$  such that for all  $\beta_1$

$$(18.2.1) \quad \text{MSE}[\tilde{\phi}; \phi | \beta = \beta_0] \geq \text{MSE}[\hat{\phi}; \phi | \beta = \beta_1]$$

In other words, the worst that can happen if one uses any other  $\tilde{\phi}$  is worse than the worst that can happen if one uses  $\hat{\phi}$ . Using this concept one can prove the following:

**THEOREM 18.2.2.**  $\hat{\beta}$  is a linear minimax estimator of the parameter vector  $\beta$  in the following sense: for every nonrandom coefficient vector  $\mathbf{t}$ ,  $\mathbf{t}^\top \hat{\beta}$  is the linear minimax estimator of the scalar  $\phi = \mathbf{t}^\top \beta$  with respect to the MSE. I.e., for every other linear estimator  $\tilde{\phi} = \mathbf{a}^\top \mathbf{y}$  of  $\phi$  one can find a value  $\beta = \beta_0$  for which  $\tilde{\phi}$  has a larger MSE than the largest possible MSE of  $\mathbf{t}^\top \hat{\beta}$ .

Proof: as in the proof of Theorem 18.1.1, write the alternative linear estimator  $\tilde{\phi}$  in the form  $\tilde{\phi} = (\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) \mathbf{y}$ , so that the sampling error is given by (18.1.5). Then it follows

$$(18.2.2) \quad \begin{aligned} \text{MSE}[\tilde{\phi}; \phi] &= \text{E}[(\tilde{\phi} - \phi)^2] = \text{E}\left[\left((\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) \boldsymbol{\varepsilon} + \mathbf{c}^\top \mathbf{X}\beta\right) \left(\boldsymbol{\varepsilon}^\top (\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t} + \mathbf{c}) + \mathbf{c}^\top \mathbf{X}\beta\right)\right] \\ (18.2.3) \quad &= \sigma^2 (\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{c}^\top) (\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t} + \mathbf{c}) + \mathbf{c}^\top \mathbf{X}\beta\beta^\top \mathbf{X}^\top \mathbf{c} \end{aligned}$$

Now there are two cases: if  $\mathbf{c}^\top \mathbf{X} = \mathbf{o}^\top$ , then  $\text{MSE}[\tilde{\phi}; \phi] = \sigma^2 \mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t} + \sigma^2 \mathbf{c}^\top \mathbf{c}$ . This does not depend on  $\beta$  and if  $\mathbf{c} \neq \mathbf{o}$  then this MSE is larger than that for  $\mathbf{c} = \mathbf{o}$ . If  $\mathbf{c}^\top \mathbf{X} \neq \mathbf{o}^\top$ , then  $\text{MSE}[\tilde{\phi}; \phi]$  is unbounded, i.e., for any finite number  $\omega$  one can always find a  $\beta_0$  for which  $\text{MSE}[\tilde{\phi}; \phi] > \omega$ . Since  $\text{MSE}[\hat{\phi}; \phi]$  is bounded, a  $\beta_0$  can be found that satisfies (18.2.1).

If we characterize the BLUE as a minimax estimator, we are using a consistent and unified principle. It is based on the concept of the MSE alone, not on a mixture between the concepts of unbiasedness and the MSE. This explains why the mathematical theory of the least squares estimator is so rich.

On the other hand, a minimax strategy is not a good estimation strategy. Nature is not the adversary of the researcher; it does not maliciously choose  $\beta$  in such a way that the researcher will be misled. This explains why the least squares principle, despite the beauty of its mathematical theory, does not give terribly good estimators (in fact, they are inadmissible, see the Section about the Stein rule below).

$\hat{\beta}$  is therefore simultaneously the solution to two very different minimization problems. We will refer to it as the OLS estimate if we refer to its property of minimizing the sum of squared errors, and as the BLUE estimator if we think of it as the best linear unbiased estimator.

Note that even if  $\sigma^2$  were known, one could not get a better linear unbiased estimator of  $\beta$ .

### 18.3. Miscellaneous Properties of the BLUE

PROBLEM 249.

- a. 1 point Instead of (14.2.22) one sometimes sees the formula

$$(18.3.1) \quad \hat{\beta} = \frac{\sum (x_t - \bar{x})y_t}{\sum (x_t - \bar{x})^2}.$$

for the slope parameter in the simple regression. Show that these formulas are mathematically equivalent.

ANSWER. Equivalence of (18.3.1) and (14.2.22) follows from  $\sum (x_t - \bar{x}) = 0$  and therefore also  $\bar{y} \sum (x_t - \bar{x}) = 0$ . Alternative proof, using matrix notation and the matrix  $D$  defined in Problem 161: (14.2.22) is  $\frac{\mathbf{x}^\top D^\top D \mathbf{y}}{\mathbf{x}^\top D^\top D \mathbf{x}}$  and (18.3.1) is  $\frac{\mathbf{x}^\top D \mathbf{y}}{\mathbf{x}^\top D \mathbf{x}}$ . They are equal because  $D$  is symmetric and idempotent. □

- b. 1 point Show that

$$(18.3.2) \quad \text{var}[\hat{\beta}] = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

ANSWER. Write (18.3.1) as

$$(18.3.3) \quad \hat{\beta} = \frac{1}{\sum (x_t - \bar{x})^2} \sum (x_t - \bar{x})y_t \quad \Rightarrow \quad \text{var}[\hat{\beta}] = \frac{1}{(\sum (x_t - \bar{x})^2)^2} \sum (x_t - \bar{x})^2 \sigma^2$$
□

- c. 2 points Show that  $\text{cov}[\hat{\beta}, \bar{y}] = 0$ .

ANSWER. This is a special case of problem 247, but it can be easily shown here separately

$$\begin{aligned} \text{cov}[\hat{\beta}, \bar{y}] &= \text{cov} \left[ \frac{\sum_s (x_s - \bar{x})y_s}{\sum_t (x_t - \bar{x})^2}, \frac{1}{n} \sum_j y_j \right] = \frac{1}{n \sum_t (x_t - \bar{x})^2} \text{cov} \left[ \sum_s (x_s - \bar{x})y_s, \sum_j y_j \right] = \\ &= \frac{1}{n \sum_t (x_t - \bar{x})^2} \sum_s (x_s - \bar{x}) \sigma^2 = 0. \end{aligned}$$

- d. 2 points Using (14.2.23) show that

$$(18.3.4) \quad \text{var}[\hat{\alpha}] = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)$$

PROBLEM 250. You have two data vectors  $x_i$  and  $y_i$  ( $i = 1, \dots, n$ ), and the following model is

$$(18.3.5) \quad y_i = \beta x_i + \varepsilon_i$$

where  $x_i$  and  $\varepsilon_i$  satisfy the basic assumptions of the linear regression model. The least squares estimator for this model is

$$(18.3.6) \quad \tilde{\beta} = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y} = \frac{\sum x_i y_i}{\sum x_i^2}$$

- a. 1 point Is  $\tilde{\beta}$  an unbiased estimator of  $\beta$ ? (Proof is required.)

ANSWER. First derive a nice expression for  $\tilde{\beta} - \beta$ :

$$\begin{aligned} \tilde{\beta} - \beta &= \frac{\sum x_i y_i}{\sum x_i^2} - \frac{\sum x_i^2 \beta}{\sum x_i^2} \\ &= \frac{\sum x_i (y_i - x_i \beta)}{\sum x_i^2} \\ &= \frac{\sum x_i \varepsilon_i}{\sum x_i^2} \quad \text{since } y_i = \beta x_i + \varepsilon_i \end{aligned}$$

$$\begin{aligned} E[\tilde{\beta} - \beta] &= E \left[ \frac{\sum x_i \varepsilon_i}{\sum x_i^2} \right] \\ &= \frac{\sum E[x_i \varepsilon_i]}{\sum x_i^2} \\ &= \frac{\sum x_i E[\varepsilon_i]}{\sum x_i^2} = 0 \quad \text{since } E \varepsilon_i = 0. \end{aligned}$$

- b. 2 points Derive the variance of  $\tilde{\beta}$ . (Show your work.)

ANSWER.

$$\begin{aligned} \text{var } \tilde{\beta} &= E[\tilde{\beta} - \beta]^2 \\ &= E\left(\frac{\sum x_i \varepsilon_i}{\sum x_i^2}\right)^2 \\ &= \frac{1}{(\sum x_i^2)^2} E\left[\sum x_i \varepsilon_i\right]^2 \\ &= \frac{1}{(\sum x_i^2)^2} \left( E\sum (x_i \varepsilon_i)^2 + 2 E\sum_{i < j} (x_i \varepsilon_i)(x_j \varepsilon_j) \right) \\ &= \frac{1}{(\sum x_i^2)^2} \sum E[x_i \varepsilon_i]^2 \quad \text{since the } \varepsilon_i \text{'s are uncorrelated, i.e., } \text{cov}[\varepsilon_i, \varepsilon_j] = 0 \text{ for } i \neq j \\ &= \frac{1}{(\sum x_i^2)^2} \sigma^2 \sum x_i^2 \quad \text{since all } \varepsilon_i \text{ have equal variance } \sigma^2 \\ &= \frac{\sigma^2}{\sum x_i^2}. \end{aligned}$$

□

PROBLEM 251. We still assume (18.3.5) is the true model. Consider an alternative estimator:

$$(18.3.7) \quad \hat{\beta} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

i.e., the estimator which would be the best linear unbiased estimator if the true model were (14.2.15).

• a. 2 points Is  $\hat{\beta}$  still an unbiased estimator of  $\beta$  if (18.3.5) is the true model? (A short but rigorous argument may save you a lot of algebra here).

ANSWER. One can argue it:  $\hat{\beta}$  is unbiased for model (14.2.15) whatever the value of  $\alpha$  or therefore also when  $\alpha = 0$ , i.e., when the model is (18.3.5). But here is the pedestrian way:

$$\begin{aligned} \hat{\beta} &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})y_i}{\sum (x_i - \bar{x})^2} \quad \text{since } \sum (x_i - \bar{x})\bar{y} = 0 \\ &= \frac{\sum (x_i - \bar{x})(\beta x_i + \varepsilon_i)}{\sum (x_i - \bar{x})^2} \quad \text{since } y_i = \beta x_i + \varepsilon_i \\ &= \beta \frac{\sum (x_i - \bar{x})x_i}{\sum (x_i - \bar{x})^2} + \frac{\sum (x_i - \bar{x})\varepsilon_i}{\sum (x_i - \bar{x})^2} \\ &= \beta + \frac{\sum (x_i - \bar{x})\varepsilon_i}{\sum (x_i - \bar{x})^2} \quad \text{since } \sum (x_i - \bar{x})x_i = \sum (x_i - \bar{x})^2 \\ E \hat{\beta} &= E\beta + E \frac{\sum (x_i - \bar{x})\varepsilon_i}{\sum (x_i - \bar{x})^2} \\ &= \beta + \frac{\sum (x_i - \bar{x}) E \varepsilon_i}{\sum (x_i - \bar{x})^2} = \beta \quad \text{since } E \varepsilon_i = 0 \text{ for all } i, \text{ i.e., } \hat{\beta} \text{ is unbiased.} \end{aligned}$$

• b. 2 points Derive the variance of  $\hat{\beta}$  if (18.3.5) is the true model.

ANSWER. One can again argue it: since the formula for  $\text{var } \hat{\beta}$  does not depend on what true value of  $\alpha$  is, it is the same formula.

$$(18.3.8) \quad \text{var } \hat{\beta} = \text{var} \left( \beta + \frac{\sum (x_i - \bar{x})\varepsilon_i}{\sum (x_i - \bar{x})^2} \right)$$

$$(18.3.9) \quad = \text{var} \left( \frac{\sum (x_i - \bar{x})\varepsilon_i}{\sum (x_i - \bar{x})^2} \right)$$

$$(18.3.10) \quad = \frac{\sum (x_i - \bar{x})^2 \text{var } \varepsilon_i}{(\sum (x_i - \bar{x})^2)^2} \quad \text{since } \text{cov}[\varepsilon_i, \varepsilon_j] = 0$$

$$(18.3.11) \quad = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}.$$

• c. 1 point Still assuming (18.3.5) is the true model, would you prefer  $\hat{\beta}$  or  $\tilde{\beta}$  from Problem 250 as an estimator of  $\beta$ ?

ANSWER. Since  $\tilde{\beta}$  and  $\hat{\beta}$  are both unbiased estimators, if (18.3.5) is the true model, the preferred estimator is the one with the smaller variance. As I will show,  $\text{var } \tilde{\beta} \leq \text{var } \hat{\beta}$  and, therefore,  $\tilde{\beta}$  is preferred to  $\hat{\beta}$ . To show

$$(18.3.12) \quad \text{var } \hat{\beta} = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \geq \frac{\sigma^2}{\sum x_i^2} = \text{var } \tilde{\beta}$$

one must show

$$(18.3.13) \quad \sum (x_i - \bar{x})^2 \leq \sum x_i^2$$

which is a simple consequence of (9.1.1). Thus  $\text{var } \hat{\beta} \geq \text{var } \tilde{\beta}$ ; the variances are equal only if  $\bar{x} = 0$ , i.e., if  $\tilde{\beta} = \hat{\beta}$ .  $\square$

**PROBLEM 252.** Suppose the true model is (14.2.15) and the basic assumptions are satisfied.

• a. 2 points In this situation,  $\tilde{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$  is generally a biased estimator of  $\beta$ . Show that its bias is

$$(18.3.14) \quad \mathbb{E}[\tilde{\beta} - \beta] = \alpha \frac{n\bar{x}}{\sum x_i^2}$$

**ANSWER.** In situations like this it is always worth while to get a nice simple expression for the sampling error:

$$(18.3.15) \quad \tilde{\beta} - \beta = \frac{\sum x_i y_i}{\sum x_i^2} - \beta$$

$$(18.3.16) \quad = \frac{\sum x_i(\alpha + \beta x_i + \varepsilon_i)}{\sum x_i^2} - \beta \quad \text{since } y_i = \alpha + \beta x_i + \varepsilon_i$$

$$(18.3.17) \quad = \alpha \frac{\sum x_i}{\sum x_i^2} + \beta \frac{\sum x_i^2}{\sum x_i^2} + \frac{\sum x_i \varepsilon_i}{\sum x_i^2} - \beta$$

$$(18.3.18) \quad = \alpha \frac{\sum x_i}{\sum x_i^2} + \frac{\sum x_i \varepsilon_i}{\sum x_i^2}$$

$$(18.3.19) \quad \mathbb{E}[\tilde{\beta} - \beta] = \mathbb{E} \alpha \frac{\sum x_i}{\sum x_i^2} + \mathbb{E} \frac{\sum x_i \varepsilon_i}{\sum x_i^2}$$

$$(18.3.20) \quad = \alpha \frac{\sum x_i}{\sum x_i^2} + \frac{\sum x_i \mathbb{E} \varepsilon_i}{\sum x_i^2}$$

$$(18.3.21) \quad = \alpha \frac{\sum x_i}{\sum x_i^2} + 0 = \alpha \frac{n\bar{x}}{\sum x_i^2}$$

This is  $\neq 0$  unless  $\bar{x} = 0$  or  $\alpha = 0$ .  $\square$

• b. 2 points Compute  $\text{var}[\tilde{\beta}]$ . Is it greater or smaller than

$$(18.3.22) \quad \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

which is the variance of the OLS estimator in this model?

**ANSWER.**

$$(18.3.23) \quad \text{var } \tilde{\beta} = \text{var} \left[ \frac{\sum x_i y_i}{\sum x_i^2} \right]$$

$$(18.3.24) \quad = \frac{1}{(\sum x_i^2)^2} \text{var} \left[ \sum x_i y_i \right]$$

$$(18.3.25) \quad = \frac{1}{(\sum x_i^2)^2} \sum x_i^2 \text{var}[y_i]$$

$$(18.3.26) \quad = \frac{\sigma^2}{(\sum x_i^2)^2} \sum x_i^2 \quad \text{since all } y_i \text{ are uncorrelated and have equal variance } \sigma^2$$

$$(18.3.27) \quad = \frac{\sigma^2}{\sum x_i^2}.$$

This variance is smaller or equal because  $\sum x_i^2 \geq \sum (x_i - \bar{x})^2$ .

• c. 5 points Show that the MSE of  $\tilde{\beta}$  is smaller than that of the OLS estimator if and only if the unknown true parameters  $\alpha$  and  $\sigma^2$  satisfy the equation

$$(18.3.28) \quad \frac{\alpha^2}{\sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)} < 1$$

**ANSWER.** This implies some tedious algebra. Here it is important to set it up right.

$$\begin{aligned} \text{MSE}[\tilde{\beta}; \beta] &= \frac{\sigma^2}{\sum x_i^2} + \left( \frac{\alpha n \bar{x}}{\sum x_i^2} \right)^2 \leq \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \\ &\left( \frac{\alpha n \bar{x}}{\sum x_i^2} \right)^2 \leq \frac{\sigma^2}{\sum (x_i - \bar{x})^2} - \frac{\sigma^2}{\sum x_i^2} = \frac{\sigma^2 (\sum x_i^2 - \sum (x_i - \bar{x})^2)}{\sum (x_i - \bar{x})^2 \sum x_i^2} \\ &= \frac{\sigma^2 n \bar{x}^2}{\sum (x_i - \bar{x})^2 \sum x_i^2} \\ \frac{\alpha^2 n}{\sum x_i^2} &= \frac{\alpha^2}{\frac{1}{n} \sum (x_i - \bar{x})^2 + \bar{x}^2} \leq \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \\ &\frac{\alpha^2}{\sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)} \leq 1 \end{aligned}$$

Now look at this lefthand side; it is amazing and surprising that it is exactly the population equivalent of the F-test for testing  $\alpha = 0$  in the regression *with* intercept. It can be estimated replacing  $\alpha^2$  with  $\hat{\alpha}^2$  and  $\sigma^2$  with  $s^2$  (in the regression with intercept). Let's look at this statistic. If  $\alpha = 0$  it has a F-distribution with 1 and  $n - 2$  degrees of freedom. If  $\alpha \neq 0$  it has what is called a noncentral distribution, and the only thing we needed to know so far was that it was likely to assume larger values than with  $\alpha = 0$ . This is why a small value of that statistic supported the hypothesis that  $\alpha = 0$ . But in the present case we are not testing whether  $\alpha = 0$  but whether the constrained MSE is better than the unconstrained. This is the case of the above inequality holds, the limiting case being that it is an equality. If it is an equality, then the above statistic has a distribution with noncentrality parameter 1/2. (Here all we need to know that: if  $z \sim N(\mu, 1)$  then

$z^2 \sim \chi_1^2$  with noncentrality parameter  $\mu^2/2$ . A noncentral F has a noncentral  $\chi^2$  in numerator and a central one in denominator.) The testing principle is therefore: compare the observed value with the upper  $\alpha$  point of a F distribution with noncentrality parameter  $1/2$ . This gives higher critical values than testing for  $\alpha = 0$ ; i.e., one may reject that  $\alpha = 0$  but not reject that the MSE of the constrained estimator is larger. This is as it should be. Compare [Gre97, 8.5.1 pp. 405–408] on this.  $\square$

From the Gauss-Markov theorem follows that for every nonrandom matrix  $\mathbf{R}$ , the BLUE of  $\boldsymbol{\phi} = \mathbf{R}\boldsymbol{\beta}$  is  $\hat{\boldsymbol{\phi}} = \mathbf{R}\hat{\boldsymbol{\beta}}$ . Furthermore, the best linear unbiased predictor (BLUP) of  $\boldsymbol{\varepsilon} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$  is the vector of residuals  $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ .

PROBLEM 253. Let  $\tilde{\boldsymbol{\varepsilon}} = \mathbf{A}\mathbf{y}$  be a linear predictor of the disturbance vector  $\boldsymbol{\varepsilon}$  in the model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2\mathbf{I})$ .

• a. 2 points Show that  $\tilde{\boldsymbol{\varepsilon}}$  is unbiased, i.e.,  $E[\tilde{\boldsymbol{\varepsilon}} - \boldsymbol{\varepsilon}] = \mathbf{o}$ , regardless of the value of  $\boldsymbol{\beta}$ , if and only if  $\mathbf{A}$  satisfies  $\mathbf{A}\mathbf{X} = \mathbf{O}$ .

ANSWER.  $E[\mathbf{A}\mathbf{y} - \boldsymbol{\varepsilon}] = E[\mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{A}\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}] = \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{o} - \mathbf{o}$ . This is  $\mathbf{o}$  for all  $\boldsymbol{\beta}$  if and only if  $\mathbf{A}\mathbf{X} = \mathbf{O}$   $\square$

• b. 2 points Which unbiased linear predictor  $\tilde{\boldsymbol{\varepsilon}} = \mathbf{A}\mathbf{y}$  of  $\boldsymbol{\varepsilon}$  minimizes the MSE-matrix  $E[(\tilde{\boldsymbol{\varepsilon}} - \boldsymbol{\varepsilon})(\tilde{\boldsymbol{\varepsilon}} - \boldsymbol{\varepsilon})^\top]$ ? Hint: Write  $\mathbf{A} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{C}$ . What is the minimum value of this MSE-matrix?

ANSWER. Since  $\mathbf{A}\mathbf{X} = \mathbf{O}$ , the prediction error  $\mathbf{A}\mathbf{y} - \boldsymbol{\varepsilon} = \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{A}\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon} = (\mathbf{A} - \mathbf{I})\boldsymbol{\varepsilon}$ ; therefore one minimizes  $\sigma^2(\mathbf{A} - \mathbf{I})(\mathbf{A} - \mathbf{I})^\top$  s. t.  $\mathbf{A}\mathbf{X} = \mathbf{O}$ . Using the hint,  $\mathbf{C}$  must also satisfy  $\mathbf{C}\mathbf{X} = \mathbf{O}$ , and  $(\mathbf{A} - \mathbf{I})(\mathbf{A} - \mathbf{I})^\top = (\mathbf{C} - \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top)(\mathbf{C}^\top - \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top) = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{C}\mathbf{C}^\top$ , therefore one must set  $\mathbf{C} = \mathbf{O}$ . Minimum value is  $\sigma^2\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top$ .  $\square$

• c. How does this best predictor relate to the OLS estimator  $\hat{\boldsymbol{\beta}}$ ?

ANSWER. It is equal to the residual vector  $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ .  $\square$

PROBLEM 254. This is a vector generalization of problem 170. Let  $\hat{\boldsymbol{\beta}}$  the BLUE of  $\boldsymbol{\beta}$  and  $\tilde{\boldsymbol{\beta}}$  an arbitrary linear unbiased estimator of  $\boldsymbol{\beta}$ .

• a. 2 points Show that  $\mathcal{C}[\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}, \hat{\boldsymbol{\beta}}] = \mathbf{O}$ .

ANSWER. Say  $\tilde{\boldsymbol{\beta}} = \tilde{\mathbf{B}}\mathbf{y}$ ; unbiasedness means  $\tilde{\mathbf{B}}\mathbf{X} = \mathbf{I}$ . Therefore

$$\begin{aligned} \mathcal{C}[\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}, \hat{\boldsymbol{\beta}}] &= \mathcal{C}[(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top - \tilde{\mathbf{B}}]\mathbf{y}, (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}] \\ &= ((\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top - \tilde{\mathbf{B}}) \mathcal{V}[\mathbf{y}]\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1} \\ &= \sigma^2((\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top - \tilde{\mathbf{B}})\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1} \\ &= \sigma^2((\mathbf{X}^\top\mathbf{X})^{-1} - (\mathbf{X}^\top\mathbf{X})^{-1}) = \mathbf{O}. \end{aligned}$$

$\square$

• b. 2 points Show that  $\mathcal{MSE}[\tilde{\boldsymbol{\beta}}; \boldsymbol{\beta}] = \mathcal{MSE}[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}] + \mathcal{V}[\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}]$

ANSWER. Due to unbiasedness,  $\mathcal{MSE} = \mathcal{V}$ , and the decomposition  $\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})$  is uncorrelated sum. Here is more detail:  $\mathcal{MSE}[\tilde{\boldsymbol{\beta}}; \boldsymbol{\beta}] = \mathcal{V}[\tilde{\boldsymbol{\beta}}] = \mathcal{V}[\hat{\boldsymbol{\beta}} + \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}] = \mathcal{V}[\hat{\boldsymbol{\beta}}] + \mathcal{C}[\tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}] + \mathcal{V}[\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}]$  but the two  $\mathcal{C}$ -terms are the null matrices.

PROBLEM 255. 3 points Given a simple regression  $y_t = \alpha + \beta x_t + \varepsilon_t$ , where  $\varepsilon_t$  are independent and identically distributed with mean  $\mu$  and variance  $\sigma^2$ . Is it possible to consistently estimate all four parameters  $\alpha, \beta, \sigma^2$ , and  $\mu$ ? If yes, explain how you would estimate them, and if no, what is the best you can do?

ANSWER. Call  $\tilde{\varepsilon}_t = \varepsilon_t - \mu$ , then the equation reads  $y_t = \alpha + \mu + \beta x_t + \tilde{\varepsilon}_t$ , with well behaved disturbances. Therefore one can estimate  $\alpha + \mu, \beta$ , and  $\sigma^2$ . This is also the best one can do. If  $\alpha + \mu$  are equal, the  $y_t$  have the same joint distribution.

PROBLEM 256. 3 points The model is  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  but all rows of the  $\mathbf{X}$ -matrix are exactly equal. What can you do? Can you estimate  $\boldsymbol{\beta}$ ? If not, are there any linear combinations of the components of  $\boldsymbol{\beta}$  which you can estimate? Can you estimate  $\sigma^2$ ?

ANSWER. If all rows are equal, then each column is a multiple of  $\mathbf{1}$ . Therefore, if there are more than one column, none of the individual components of  $\boldsymbol{\beta}$  can be estimated. But you can estimate  $\mathbf{x}^\top\boldsymbol{\beta}$  (if  $\mathbf{x}$  is one of the row vectors of  $\mathbf{X}$ ) and you can estimate  $\sigma^2$ .

PROBLEM 257. This is [JHG+88, 5.3.32]: Consider the log-linear statistical model

$$(18.3.29) \quad \mathbf{y}_t = \alpha x_t^\beta \exp \varepsilon_t = z_t \exp \varepsilon_t$$

with “well-behaved” disturbances  $\varepsilon_t$ . Here  $z_t = \alpha x_t^\beta$  is the systematic portion of  $y_t$  which depends on  $x_t$ . (This functional form is often used in models of demand and production.)

• a. 1 point Can this be estimated with the regression formalism?

ANSWER. Yes, simply take logs:

$$(18.3.30) \quad \log y_t = \log \alpha + \beta \log x_t + \varepsilon_t$$

• b. 1 point Show that the elasticity of the functional relationship between  $x_t$  and  $z_t$

$$(18.3.31) \quad \eta = \frac{\partial z_t / z_t}{\partial x_t / x_t}$$

does not depend on  $t$ , i.e., it is the same for all observations. Many authors talk about the elasticity of  $y_t$  with respect to  $x_t$ , but one should really only talk about elasticity of  $z_t$  with respect to  $x_t$ , where  $z_t$  is the systematic part of  $y_t$  which can be estimated by  $\hat{y}_t$ .

ANSWER. The systematic functional relationship is  $\log z_t = \log \alpha + \beta \log x_t$ ; therefore

$$(18.3.32) \quad \frac{\partial \log z_t}{\partial z_t} = \frac{1}{z_t}$$

which can be rewritten as

$$(18.3.33) \quad \frac{\partial z_t}{z_t} = \partial \log z_t;$$

The same can be done with  $x_t$ ; therefore

$$(18.3.34) \quad \frac{\partial z_t / z_t}{\partial x_t / x_t} = \frac{\partial \log z_t}{\partial \log x_t} = \beta$$

What we just did was a tricky way to take a derivative. A less tricky way is:

$$(18.3.35) \quad \frac{\partial z_t}{\partial x_t} = \alpha \beta x_t^{\beta-1} = \beta z_t / x_t$$

Therefore

$$(18.3.36) \quad \frac{\partial z_t}{\partial x_t} \frac{x_t}{z_t} = \beta$$

□

PROBLEM 258.

- a. 2 points What is the elasticity in the simple regression  $y_t = \alpha + \beta x_t + \varepsilon_t$ ?

ANSWER.

$$(18.3.37) \quad \eta_t = \frac{\partial z_t / z_t}{\partial x_t / x_t} = \frac{\partial z_t}{\partial x_t} \frac{x_t}{z_t} = \frac{\beta x_t}{z_t} = \frac{\beta x_t}{\alpha + \beta x_t}$$

This depends on the observation, and if one wants one number, a good way is to evaluate it at  $\bar{x}$ . □

- b. Show that an estimate of this elasticity evaluated at  $\bar{x}$  is  $h = \frac{\hat{\beta}\bar{x}}{\bar{y}}$ .

ANSWER. This comes from the fact that the fitted regression line goes through the point  $\bar{x}, \bar{y}$ . If one uses the other definition of elasticity, which Greene uses on p. 227 but no longer on p. 280, and which I think does not make much sense, one gets the same formula:

$$(18.3.38) \quad \eta_t = \frac{\partial y_t / y_t}{\partial x_t / x_t} = \frac{\partial y_t}{\partial x_t} \frac{x_t}{y_t} = \frac{\beta x_t}{y_t}$$

This is different than (18.3.37), but if one evaluates it at the sample mean, both formulas give the same result  $\frac{\hat{\beta}\bar{x}}{\bar{y}}$ . □

- c. Show by the delta method that the estimator

$$(18.3.39) \quad h = \frac{\hat{\beta}\bar{x}}{\bar{y}}$$

of the elasticity in the simple regression model has the estimated asymptotic variance

$$(18.3.40) \quad s^2 \begin{bmatrix} -h & \frac{\bar{x}(1-h)}{\bar{y}} \end{bmatrix} \begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{bmatrix}^{-1} \begin{bmatrix} \frac{-h}{\bar{y}} \\ \frac{\bar{x}(1-h)}{\bar{y}} \end{bmatrix}$$

- d. Compare [Gre97, example 6.20 on p. 280]. Assume

$$(18.3.41) \quad \frac{1}{n}(\mathbf{X}^\top \mathbf{X}) = \begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{bmatrix} \rightarrow \mathbf{Q} = \begin{bmatrix} 1 & q \\ q & r \end{bmatrix}$$

where we assume for the sake of the argument that  $q$  is known. The true elasticity of the underlying functional relationship, evaluated at  $\lim \bar{x}$ , is

$$(18.3.42) \quad \eta = \frac{q\beta}{\alpha + q\beta}$$

Then

$$(18.3.43) \quad h = \frac{q\hat{\beta}}{\hat{\alpha} + q\hat{\beta}}$$

is a consistent estimate for  $\eta$ .

A generalization of the log-linear model is the translog model, which is a second order approximation to an unknown functional form, and which allows to model second-order effects such as elasticities of substitution etc. Used to model production cost, and utility functions. Start with any function  $v = f(u_1, \dots, u_n)$  and make second-order Taylor development around  $\mathbf{u} = \mathbf{o}$ :

$$(18.3.44) \quad v = f(\mathbf{o}) + \sum u_i \frac{\partial f}{\partial u_i} \Big|_{\mathbf{u}=\mathbf{o}} + \frac{1}{2} \sum_{i,j} u_i u_j \frac{\partial^2 f}{\partial u_i \partial u_j} \Big|_{\mathbf{u}=\mathbf{o}}$$

Now say  $v = \log(\mathbf{y})$  and  $u_i = \log(x_i)$ , and the values of  $f$  and its derivatives at  $\mathbf{o}$  are the coefficients to be estimated:

$$(18.3.45) \quad \log(\mathbf{y}) = \alpha + \sum \beta_i \log x_i + \frac{1}{2} \sum_{i,j} \gamma_{ij} \log x_i \log x_j + \varepsilon$$

Note that by Young's theorem it must be true that  $\gamma_{kl} = \gamma_{lk}$ .

The semi-log model is often used to model growth rates:

$$(18.3.46) \quad \log y_t = \mathbf{x}_t^\top \boldsymbol{\beta} + \varepsilon_t$$

Here usually one of the columns of  $\mathbf{X}$  is the time subscript  $t$  itself; [Gre97, p. 2] writes it as

$$(18.3.47) \quad \log y_t = \mathbf{x}_t^\top \boldsymbol{\beta} + t\delta + \varepsilon_t$$

where  $\delta$  is the autonomous growth rate. The logistic functional form is appropriate for adoption rates  $0 \leq y_t \leq 1$ : the rate of adoption is slow at first, then rapid as innovation gains popularity, then slow again as the market becomes saturated:

$$(18.3.48) \quad y_t = \frac{\exp(\mathbf{x}_t^\top \boldsymbol{\beta} + t\delta + \varepsilon_t)}{1 + \exp(\mathbf{x}_t^\top \boldsymbol{\beta} + t\delta + \varepsilon_t)}$$

This can be linearized by the logit transformation:

$$(18.3.49) \quad \text{logit}(\mathbf{y}_t) = \log \frac{\mathbf{y}_t}{1 - \mathbf{y}_t} = \mathbf{x}_t^\top \boldsymbol{\beta} + t\delta + \varepsilon_t$$

**PROBLEM 259.** 3 points Given a simple regression  $\mathbf{y}_t = \alpha_t + \beta x_t$  which deviates from an ordinary regression in two ways: (1) There is no disturbance term. (2) The “constant term”  $\alpha_t$  is random, i.e., in each time period  $t$ , the value of  $\alpha_t$  is obtained by an independent drawing from a population with unknown mean  $\mu$  and unknown variance  $\sigma^2$ . Is it possible to estimate all three parameters  $\beta$ ,  $\sigma^2$ , and  $\mu$ , and to “predict” each  $\alpha_t$ ? (Here I am using the term “prediction” for the estimation of a random parameter.) If yes, explain how you would estimate it, and if not, what is the best you can do?

**ANSWER.** Call  $\varepsilon_t = \alpha_t - \mu$ , then the equation reads  $\mathbf{y}_t = \mu + \beta x_t + \varepsilon_t$ , with well behaved disturbances. Therefore one can estimate all the unknown parameters, and predict  $\alpha_t$  by  $\hat{\mu} + \varepsilon_t$ . □

### 18.4. Estimation of the Variance

The formulas in this section use g-inverses (compare (A.3.1)) and are valid even if not all columns of  $\mathbf{X}$  are linearly independent.  $q$  is the rank of  $\mathbf{X}$ . The proofs are not any more complicated than in the case that  $\mathbf{X}$  has full rank, if one keeps in mind identity (A.3.3) and some other simple properties of g-inverses which are tacitly used at various places. Those readers who are only interested in the full-rank case should simply substitute  $(\mathbf{X}^\top \mathbf{X})^{-1}$  for  $(\mathbf{X}^\top \mathbf{X})^-$  and  $k$  for  $q$  ( $k$  is the number of columns of  $\mathbf{X}$ ).

SSE, the attained minimum value of the Least Squares objective function, is a random variable too and we will now compute its expected value. It turns out that

$$(18.4.1) \quad \mathbb{E}[\text{SSE}] = \sigma^2(n - q)$$

**PROOF.**  $\text{SSE} = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}$ , where  $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{y} = \mathbf{M}\mathbf{y}$ , with  $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top$ . From  $\mathbf{M}\mathbf{X} = \mathbf{O}$  follows  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{M}\boldsymbol{\varepsilon}$ . Since  $\mathbf{M}$  is idempotent and symmetric, it follows  $\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon}^\top \mathbf{M}\boldsymbol{\varepsilon}$ , therefore  $\mathbb{E}[\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}] = \mathbb{E}[\text{tr } \boldsymbol{\varepsilon}^\top \mathbf{M}\boldsymbol{\varepsilon}] = \mathbb{E}[\text{tr } \mathbf{M}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top] = \sigma^2 \text{tr } \mathbf{M} = \sigma^2 \text{tr}(\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top) = \sigma^2(n - \text{tr}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X}) = \sigma^2(n - q)$ . □

**PROBLEM 260.**

- a. 2 points Show that

$$(18.4.2) \quad \text{SSE} = \boldsymbol{\varepsilon}^\top \mathbf{M}\boldsymbol{\varepsilon} \quad \text{where} \quad \mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top$$

**ANSWER.**  $\text{SSE} = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}$ , where  $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{y} = \mathbf{M}\mathbf{y}$  where  $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top$ . From  $\mathbf{M}\mathbf{X} = \mathbf{O}$  follows  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{M}\boldsymbol{\varepsilon}$ . Since  $\mathbf{M}$  is idempotent and symmetric, it follows  $\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon}^\top \mathbf{M}\boldsymbol{\varepsilon}$ . □

- b. 1 point Is SSE observed? Is  $\boldsymbol{\varepsilon}$  observed? Is  $\mathbf{M}$  observed?
- c. 3 points Under the usual assumption that  $\mathbf{X}$  has full column rank, show that

$$(18.4.3) \quad \mathbb{E}[\text{SSE}] = \sigma^2(n - k)$$

**ANSWER.**  $\mathbb{E}[\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}] = \mathbb{E}[\text{tr } \boldsymbol{\varepsilon}^\top \mathbf{M}\boldsymbol{\varepsilon}] = \mathbb{E}[\text{tr } \mathbf{M}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top] = \sigma^2 \text{tr } \mathbf{M} = \sigma^2 \text{tr}(\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top) = \sigma^2(n - \text{tr}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X}) = \sigma^2(n - k)$ .

**PROBLEM 261.** As an alternative proof of (18.4.3) show that  $\text{SSE} = \mathbf{y}^\top \mathbf{M}\mathbf{y}$  and use theorem ??.

From (18.4.3) follows that  $\text{SSE}/(n - q)$  is an unbiased estimate of  $\sigma^2$ . Although it is commonly suggested that  $s^2 = \text{SSE}/(n - q)$  is an optimal estimator of  $\sigma^2$ , this is a fallacy. The question which estimator of  $\sigma^2$  is best depends on the kurtosis of the distribution of the error terms. For instance, if the kurtosis is zero, which is the case when the error terms are normal, then a different scalar multiple of the SSE, namely, the Theil-Schweitzer estimator from [TS61]

$$(18.4.4) \quad \hat{\sigma}_{TS}^2 = \frac{1}{n - q + 2} \mathbf{y}^\top \mathbf{M}\mathbf{y} = \frac{1}{n - q + 2} \sum_{i=1}^n \hat{\varepsilon}_i^2,$$

is biased but has lower MSE than  $s^2$ . Compare problem 163. The only thing one can say about  $s^2$  is that it is a fairly good estimator which one can use when one does not know the kurtosis (but even in this case it is not the best one can do).

### 18.5. Mallows’s Cp-Statistic as Estimator of the Mean Squared Error

**PROBLEM 262.** We will compute here the  $\text{MSE}$ -matrix of  $\hat{\mathbf{y}}$  as an estimator of  $\mathcal{E}[\mathbf{y}]$  in a regression which does not use the correct  $\mathbf{X}$ -matrix. For this we assume that  $\mathbf{y} = \boldsymbol{\eta} + \boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$ .  $\boldsymbol{\eta} = \mathcal{E}[\mathbf{y}]$  is an arbitrary vector of constants, and do not assume that  $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$  for some  $\boldsymbol{\beta}$ , i.e., we do not assume that  $\mathbf{X}$  contains the necessary explanatory variables. Regression of  $\mathbf{y}$  on  $\mathbf{X}$  gives the OLS estimator  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{y}$ .

- a. 2 points Show that the  $\text{MSE}$  matrix of  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  as estimator of  $\boldsymbol{\eta}$  is

$$(18.5.1) \quad \text{MSE}[\mathbf{X}\hat{\boldsymbol{\beta}}; \boldsymbol{\eta}] = \sigma^2 \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top + \mathbf{M}\boldsymbol{\eta}\boldsymbol{\eta}^\top \mathbf{M}$$

where  $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top$ .

- b. 1 point Formula (18.5.1) for the  $\text{MSE}$  matrix depends on the unknown  $\boldsymbol{\eta}$  and is therefore useless for estimation. If one cannot get an estimate of the whole  $\text{MSE}$  matrix, an often-used second best choice is its trace. Show that

$$(18.5.2) \quad \text{tr } \text{MSE}[\mathbf{X}\hat{\boldsymbol{\beta}}; \boldsymbol{\eta}] = \sigma^2 q + \boldsymbol{\eta}^\top \mathbf{M}\boldsymbol{\eta}.$$

where  $q$  is the rank of  $\mathbf{X}$ .

• c. 3 points If an unbiased estimator of the true  $\sigma^2$  is available (call it  $s^2$ ), then an unbiased estimator of the righthand side of (18.5.2) can be constructed using this  $s^2$  and the SSE of the regression  $\text{SSE} = \mathbf{y}^\top \mathbf{M} \mathbf{y}$ . Show that

$$(18.5.3) \quad E[\text{SSE} - (n - 2q)s^2] = \sigma^2 q + \boldsymbol{\eta}^\top \mathbf{M} \boldsymbol{\eta}.$$

*Hint: use equation (??). If one does not have an unbiased estimator  $s^2$  of  $\sigma^2$ , one usually gets such an estimator by regressing  $\mathbf{y}$  on an  $\mathbf{X}$  matrix which is so large that one can assume that it contains the true regressors.*

The statistic

$$(18.5.4) \quad C_p = \frac{\text{SSE}}{s^2} + 2q - n$$

is called Mallow's  $C_p$  statistic. It is a consistent estimator of  $\text{tr } \mathcal{MSE}[\mathbf{X}\hat{\boldsymbol{\beta}}; \boldsymbol{\eta}] / \sigma^2$ . If  $\mathbf{X}$  contains all necessary variables, i.e.,  $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$  for some  $\boldsymbol{\beta}$ , then (18.5.2) becomes  $\text{tr } \mathcal{MSE}[\mathbf{X}\hat{\boldsymbol{\beta}}; \boldsymbol{\eta}] = \sigma^2 q$ , i.e., in this case  $C_p$  should be close to  $q$ . Therefore the selection rule for regressions should be here to pick that regression for which the  $C_p$ -value is closest to  $q$ . (This is an explanation; nothing to prove here.)

If one therefore has several regressions and tries to decide which is the right one, it is recommended to plot  $C_p$  versus  $q$  for all regressions, and choose one for which this value is small and lies close to the diagonal. An example of this is given in problem 232.



## CHAPTER 19

## Nonspherical Positive Definite Covariance Matrix

The so-called “Generalized Least Squares” model specifies  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2\boldsymbol{\Psi})$  where  $\sigma^2$  is an unknown positive scalar, and  $\boldsymbol{\Psi}$  is a *known* positive definite matrix.

This is simply the OLS model in disguise. To see this, we need a few more facts about positive definite matrices.  $\boldsymbol{\Psi}$  is *nonnegative* definite if and only if a  $\mathbf{Q}$  exists with  $\boldsymbol{\Psi} = \mathbf{Q}\mathbf{Q}^\top$ . If  $\boldsymbol{\Psi}$  is *positive* definite, this  $\mathbf{Q}$  can be chosen square and nonsingular. Then  $\mathbf{P} = \mathbf{Q}^{-1}$  satisfies  $\mathbf{P}^\top\mathbf{P}\boldsymbol{\Psi} = \mathbf{P}^\top\mathbf{P}\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}$ , i.e.,  $\mathbf{P}^\top\mathbf{P} = \boldsymbol{\Psi}^{-1}$ , and also  $\mathbf{P}\boldsymbol{\Psi}\mathbf{P}^\top = \mathbf{P}\mathbf{Q}\mathbf{Q}^\top\mathbf{P}^\top = \mathbf{I}$ . Premultiplying the GLS model by  $\mathbf{P}$  gives therefore a model whose disturbances have a spherical covariance matrix:

$$(19.0.5) \quad \mathbf{P}\mathbf{y} = \mathbf{P}\mathbf{X}\boldsymbol{\beta} + \mathbf{P}\boldsymbol{\varepsilon} \quad \mathbf{P}\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2\mathbf{I})$$

The OLS estimate of  $\boldsymbol{\beta}$  in this transformed model is

$$(19.0.6) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top\mathbf{P}^\top\mathbf{P}\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{P}^\top\mathbf{P}\mathbf{y} = (\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{y}.$$

This  $\hat{\boldsymbol{\beta}}$  is the BLUE of  $\boldsymbol{\beta}$  in model (19.0.5), and since estimators which are linear in  $\mathbf{P}\mathbf{y}$  are also linear in  $\mathbf{y}$  and vice versa,  $\hat{\boldsymbol{\beta}}$  is also the BLUE in the original GLS model.

PROBLEM 263. 2 points Show that

$$(19.0.7) \quad \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\boldsymbol{\varepsilon}$$

and derive from this that  $\hat{\boldsymbol{\beta}}$  is unbiased and that  $\text{MSE}[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}] = \sigma^2(\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})^{-1}$ .

ANSWER. Proof of (19.0.7) is very similar to proof of (18.0.7).  $\square$

The objective function of the associated least squares problem is

$$(19.0.8) \quad \boldsymbol{\beta} = \hat{\boldsymbol{\beta}} \quad \text{minimizes} \quad (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top\boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

The normal equations are

$$(19.0.9) \quad \mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{y}$$

If  $\mathbf{X}$  has full rank, then  $\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X}$  is nonsingular, and the unique  $\hat{\boldsymbol{\beta}}$  minimizing (19.0.8) is

$$(19.0.10) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{y}$$

PROBLEM 264. [Seb77, p. 386, 5] Show that if  $\boldsymbol{\Psi}$  is positive definite and  $\mathbf{X}$  has full rank, then also  $\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X}$  is positive definite. You are allowed to use, without proof, that the inverse of a positive definite matrix is also positive definite.

ANSWER. From  $\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X}\mathbf{a} = \mathbf{o}$  follows  $\mathbf{a}^\top\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X}\mathbf{a} = 0$ , and since  $\boldsymbol{\Psi}^{-1}$  is positive definite, it follows  $\mathbf{X}\mathbf{a} = \mathbf{o}$ , and since  $\mathbf{X}$  has full column rank, this implies  $\mathbf{a} = \mathbf{o}$ .

The least squares objective function of the transformed model, which  $\hat{\boldsymbol{\beta}}$  minimizes, can be written

$$(19.0.11) \quad (\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{X}\boldsymbol{\beta})^\top(\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{X}\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top\boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}),$$

and whether one writes it in one form or the other,  $1/(n-k)$  times the minimum value of that GLS objective function is still an unbiased estimate of  $\sigma^2$ .

PROBLEM 265. Show that the minimum value of the GLS objective function can be written in the form  $\mathbf{y}^\top\mathbf{M}\mathbf{y}$  where  $\mathbf{M} = \boldsymbol{\Psi}^{-1} - \boldsymbol{\Psi}^{-1}\mathbf{X}(\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\Psi}^{-1}$ . Does  $\mathbf{M}\mathbf{X} = \mathbf{O}$  still hold? Does  $\mathbf{M}^2 = \mathbf{M}$  or a similar simple identity still hold? Show that  $\mathbf{M}$  is nonnegative definite. Show that  $\text{E}[\mathbf{y}^\top\mathbf{M}\mathbf{y}] = (n-k)\sigma^2$ .

ANSWER. In  $(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top\boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$  plug in  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{y}$  and multiply to get  $\mathbf{y}^\top\mathbf{M}\mathbf{y}$ . Yes,  $\mathbf{M}\mathbf{X} = \mathbf{O}$  holds.  $\mathbf{M}$  is no longer idempotent, but it satisfies  $\mathbf{M}\boldsymbol{\Psi}\mathbf{M} = \mathbf{M}$ . One way to show that it is and would be to use the first part of the question: for all  $\mathbf{z}$ ,  $\mathbf{z}^\top\mathbf{M}\mathbf{z} = (\mathbf{z} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top(\mathbf{z} - \mathbf{X}\hat{\boldsymbol{\beta}})$ , and another way would be to use the second part of the question:  $\mathbf{M}$  is nonnegative definite because  $\mathbf{M}\boldsymbol{\Psi}\mathbf{M} = \mathbf{M}$ . To show expected value, show first that  $\mathbf{y}^\top\mathbf{M}\mathbf{y} = \boldsymbol{\varepsilon}^\top\mathbf{M}\boldsymbol{\varepsilon}$ , and then use the tricks with the trace again.

The simplest example of Generalized Least Squares is that where  $\boldsymbol{\Psi}$  is diagonal (heteroskedastic data). In this case, the GLS objective function  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top\boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  is simply a weighted least squares, with the weights being the inverses of the diagonal elements of  $\boldsymbol{\Psi}$ . This vector of inverse diagonal elements can be specified with the optional `weights` argument in `R`, see the help-file for `lm`. Heteroskedastic data arise for instance when each data point is an average over a different number of individuals.

If one runs OLS on the original instead of the transformed model, one gets the OLS estimator, we will call it here  $\hat{\boldsymbol{\beta}}_{OLS}$ , which is still unbiased. The estimator is usually also consistent, but no longer BLUE. This not only makes it less efficient than the GLS, but one also gets the wrong results if one relies on the standard computer printouts for significance tests etc. The estimate of  $\sigma^2$  generated by this regression is now usually biased. How biased it is depends on the  $\mathbf{X}$ -matrix, but most often it seems biased upwards. The estimated standard errors in the regression printouts

not only use the wrong  $s$ , but they also insert this wrong  $s$  into the wrong formula  $\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$  instead of  $\sigma^2(\mathbf{X}\Psi^{-1}\mathbf{X})^{-1}$  for  $\mathcal{V}[\hat{\boldsymbol{\beta}}]$ .

PROBLEM 266. In the generalized least squares model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2\Psi)$ , the BLUE is

$$(19.0.12) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Psi^{-1} \mathbf{y}.$$

We will write  $\hat{\boldsymbol{\beta}}_{OLS}$  for the ordinary least squares estimator

$$(19.0.13) \quad \hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

which has different properties now since we do not assume  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2\mathbf{I})$  but  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2\Psi)$ .

- a. 1 point Is  $\hat{\boldsymbol{\beta}}_{OLS}$  unbiased?
- b. 2 points Show that, still under the assumption  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2\Psi)$ ,  $\mathcal{V}[\hat{\boldsymbol{\beta}}_{OLS}] - \mathcal{V}[\hat{\boldsymbol{\beta}}] = \mathcal{V}[\hat{\boldsymbol{\beta}}_{OLS} - \hat{\boldsymbol{\beta}}]$ . (Write down the formulas for the left hand side and the right hand side and then show by matrix algebra that they are equal.) (This is what one should expect after Problem 170.) Since due to unbiasedness the covariance matrices are the  $\mathcal{MSE}$ -matrices, this shows that  $\mathcal{MSE}[\hat{\boldsymbol{\beta}}_{OLS}; \boldsymbol{\beta}] - \mathcal{MSE}[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}]$  is nonnegative definite.

ANSWER. Verify equality of the following two expressions for the differences in  $\mathcal{MSE}$  matrices:

$$\begin{aligned} \mathcal{V}[\hat{\boldsymbol{\beta}}_{OLS}] - \mathcal{V}[\hat{\boldsymbol{\beta}}] &= \sigma^2 \left( (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \Psi \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} - (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} \right) = \\ &= \sigma^2 \left( (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Psi^{-1} \right) \Psi \left( \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} - \Psi^{-1} \mathbf{X} (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} \right) \end{aligned}$$

□

Examples of GLS models are discussed in chapters ?? and ??.

## CHAPTER 20

## Best Linear Prediction

Best Linear Prediction is the second basic building block for the linear model, in addition to the OLS model. Instead of estimating a nonrandom parameter  $\beta$  about which no prior information is available, in the present situation one predicts a random variable  $\mathbf{z}$  whose mean and covariance matrix are known. Most models to be discussed below are somewhere between these two extremes.

Christensen's [Chr87] is one of the few textbooks which treat best linear prediction on the basis of known first and second moments in parallel with the regression model. The two models have indeed so much in common that they should be treated together.

## 20.1. Minimum Mean Squared Error, Unbiasedness Not Required

Assume the expected values of the random vectors  $\mathbf{y}$  and  $\mathbf{z}$  are known, and their joint covariance matrix is known up to an unknown scalar factor  $\sigma^2 > 0$ . We will write this as

$$(20.1.1) \quad \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\nu} \end{bmatrix}, \sigma^2 \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} \\ \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}} \end{bmatrix}, \quad \sigma^2 > 0.$$

$\mathbf{y}$  is observed but  $\mathbf{z}$  is not, and the goal is to predict  $\mathbf{z}$  on the basis of the observation of  $\mathbf{y}$ .

There is a unique predictor of the form  $\mathbf{z}^* = \mathbf{B}^*\mathbf{y} + \mathbf{b}^*$  (i.e., it is linear with a constant term, the technical term for this is "affine") with the following two properties: it is unbiased, and the prediction error is uncorrelated with  $\mathbf{y}$ , i.e.,

$$(20.1.2) \quad \mathcal{C}[\mathbf{z}^* - \mathbf{z}, \mathbf{y}] = \mathbf{O}.$$

The formulas for  $\mathbf{B}^*$  and  $\mathbf{b}^*$  are easily derived. Unbiasedness means  $\boldsymbol{\nu} = \mathbf{B}^*\boldsymbol{\mu} + \mathbf{b}^*$ , the predictor has therefore the form

$$(20.1.3) \quad \mathbf{z}^* = \boldsymbol{\nu} + \mathbf{B}^*(\mathbf{y} - \boldsymbol{\mu}).$$

Since

$$(20.1.4) \quad \mathbf{z}^* - \mathbf{z} = \mathbf{B}^*(\mathbf{y} - \boldsymbol{\mu}) - (\mathbf{z} - \boldsymbol{\nu}) = [\mathbf{B}^* \quad -\mathbf{I}] \begin{bmatrix} \mathbf{y} - \boldsymbol{\mu} \\ \mathbf{z} - \boldsymbol{\nu} \end{bmatrix},$$

the zero correlation condition (20.1.2) translates into

$$(20.1.5) \quad \mathbf{B}^*\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} = \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}},$$

which, due to equation (A.5.13) holds for  $\mathbf{B}^* = \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1}$ . Therefore the predictor

$$(20.1.6) \quad \mathbf{z}^* = \boldsymbol{\nu} + \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$$

satisfies the two requirements.

Unbiasedness and condition (20.1.2) are sometimes interpreted to mean that is an optimal predictor. Unbiasedness is often naively (but erroneously) considered to be a necessary condition for good estimators. And if the prediction error were correlated with the observed variable, the argument goes, then it would be possible to improve the prediction. Theorem 20.1.1 shows that despite the flaws in the argument the result which it purports to show is indeed valid:  $\mathbf{z}^*$  has the minimum  $\mathcal{MSE}$  among all affine predictors, whether biased or not, of  $\mathbf{z}$  on the basis of  $\mathbf{y}$ .

**THEOREM 20.1.1.** *In situation (20.1.1), the predictor (20.1.6) has, among predictors of  $\mathbf{z}$  which are affine functions of  $\mathbf{y}$ , the smallest  $\mathcal{MSE}$  matrix. Its  $\mathcal{MSE}$  matrix is*

$$(20.1.7) \quad \mathcal{MSE}[\mathbf{z}^*; \mathbf{z}] = \mathcal{E}[(\mathbf{z}^* - \mathbf{z})(\mathbf{z}^* - \mathbf{z})^\top] = \sigma^2(\boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}} - \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}) = \sigma^2\boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}|\mathbf{y}}$$

**PROOF.** Look at any predictor of the form  $\tilde{\mathbf{z}} = \tilde{\mathbf{B}}\mathbf{y} + \tilde{\mathbf{b}}$ . Its bias is  $\tilde{\mathbf{d}} = \mathcal{E}[\tilde{\mathbf{z}} - \mathbf{z}] = \tilde{\mathbf{B}}\boldsymbol{\mu} + \tilde{\mathbf{b}} - \boldsymbol{\nu}$ , and by (17.1.2) one can write

$$(20.1.8) \quad \mathcal{E}[(\tilde{\mathbf{z}} - \mathbf{z})(\tilde{\mathbf{z}} - \mathbf{z})^\top] = \nu[(\tilde{\mathbf{z}} - \mathbf{z})] + \tilde{\mathbf{d}}\tilde{\mathbf{d}}^\top$$

$$(20.1.9) \quad = \nu \left[ \begin{bmatrix} \tilde{\mathbf{B}} & -\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} \right] + \tilde{\mathbf{d}}\tilde{\mathbf{d}}^\top$$

$$(20.1.10) \quad = \sigma^2 \begin{bmatrix} \tilde{\mathbf{B}} & -\mathbf{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} \\ \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{B}}^\top \\ -\mathbf{I} \end{bmatrix} + \tilde{\mathbf{d}}\tilde{\mathbf{d}}^\top.$$

This  $\mathcal{MSE}$ -matrix is minimized if and only if  $\tilde{\mathbf{d}} = \mathbf{0}$  and  $\tilde{\mathbf{B}}^*$  satisfies (20.1.5). To do this, take any solution  $\mathbf{B}^*$  of (20.1.5), and write  $\tilde{\mathbf{B}} = \mathbf{B}^* + \tilde{\mathbf{D}}$ . Since, due to theorem A.5.11,  $\boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}} = \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}$ , it follows  $\boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\mathbf{B}^{*\top} = \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}\mathbf{B}^{*\top} = \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}$ . Therefore

$$(20.1.11) \quad \mathcal{MSE}[\tilde{\mathbf{z}}; \mathbf{z}] = \sigma^2 \begin{bmatrix} \mathbf{B}^* + \tilde{\mathbf{D}} & -\mathbf{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} \\ \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}} \end{bmatrix} \begin{bmatrix} \mathbf{B}^{*\top} + \tilde{\mathbf{D}}^\top \\ -\mathbf{I} \end{bmatrix} + \tilde{\mathbf{d}}\tilde{\mathbf{d}}^\top$$

$$(20.1.11) \quad = \sigma^2 \begin{bmatrix} \mathbf{B}^* + \tilde{\mathbf{D}} & -\mathbf{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}\tilde{\mathbf{D}}^\top \\ -\boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}|\mathbf{y}} + \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}}\tilde{\mathbf{D}}^\top \end{bmatrix} + \tilde{\mathbf{d}}\tilde{\mathbf{d}}^\top$$

$$(20.1.12) \quad = \sigma^2(\boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}|\mathbf{y}} + \tilde{\mathbf{D}}\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}\tilde{\mathbf{D}}^\top) + \tilde{\mathbf{d}}\tilde{\mathbf{d}}^\top.$$

The  $\mathcal{MSE}$  matrix is therefore minimized (with minimum value  $\sigma^2\Omega_{zz,y}$ ) if and only if  $\tilde{\mathbf{d}} = \mathbf{o}$  and  $\tilde{\mathbf{D}}\Omega_{yy} = \mathbf{O}$  which means that  $\tilde{\mathbf{B}}$ , along with  $\mathbf{B}^*$ , satisfies (20.1.5).  $\square$

PROBLEM 267. Show that the solution of this minimum MSE problem is unique in the following sense: if  $\mathbf{B}_1^*$  and  $\mathbf{B}_2^*$  are two different solutions of (20.1.5) and  $\mathbf{y}$  is any feasible observed value  $\mathbf{y}$ , plugged into equations (20.1.3) they will lead to the same predicted value  $\mathbf{z}^*$ .

ANSWER. Comes from the fact that every feasible observed value of  $\mathbf{y}$  can be written in the form  $\mathbf{y} = \boldsymbol{\mu} + \Omega_{yy}\mathbf{q}$  for some  $\mathbf{q}$ , therefore  $\mathbf{B}_i^*\mathbf{y} = \mathbf{B}_i^*\Omega_{yy}\mathbf{q} = \Omega_{zy}\mathbf{q}$ .  $\square$

The matrix  $\mathbf{B}^*$  is also called the regression matrix of  $\mathbf{z}$  on  $\mathbf{y}$ , and the unscaled covariance matrix has the form

$$(20.1.13) \quad \Omega = \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix} = \begin{bmatrix} \Omega_{yy} & \Omega_{yy}\mathbf{X}^\top \\ \mathbf{X}\Omega_{yy} & \mathbf{X}\Omega_{yy}\mathbf{X}^\top + \Omega_{zz,y} \end{bmatrix}$$

Where we wrote here  $\mathbf{B}^* = \mathbf{X}$  in order to make the analogy with regression clearer. A g-inverse is

$$(20.1.14) \quad \Omega^- = \begin{bmatrix} \Omega_{yy}^- + \mathbf{X}^\top\Omega_{zz,y}^- \mathbf{X} & -\mathbf{X}^\top\Omega_{zz,y}^- \\ -\mathbf{X}^\top\Omega_{zz,y}^- & \Omega_{zz,y}^- \end{bmatrix}$$

and every g-inverse of the covariance matrix has a g-inverse of  $\Omega_{zz,y}$  as its  $\mathbf{z}\mathbf{z}$ -partition. (Proof in Problem 392.)

If  $\Omega = \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix}$  is nonsingular, 20.1.5 is also solved by  $\mathbf{B}^* = -(\Omega^{zz})^{-1}\Omega^{zy}$

where  $\Omega^{zz}$  and  $\Omega^{zy}$  are the corresponding partitions of the inverse  $\Omega^{-1}$ . See Problem 392 for a proof. Therefore instead of 20.1.6 the predictor can also be written

$$(20.1.15) \quad \mathbf{z}^* = \boldsymbol{\nu} - (\Omega^{zz})^{-1}\Omega^{zy}(\mathbf{y} - \boldsymbol{\mu})$$

(note the minus sign) or

$$(20.1.16) \quad \mathbf{z}^* = \boldsymbol{\nu} - \Omega_{zz,y}\Omega^{zy}(\mathbf{y} - \boldsymbol{\mu}).$$

PROBLEM 268. This problem utilizes the concept of a bounded risk estimator, which is not yet explained very well in these notes. Assume  $\mathbf{y}$ ,  $\mathbf{z}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\nu}$  are jointly distributed random vectors. First assume  $\boldsymbol{\nu}$  and  $\boldsymbol{\mu}$  are observed, but  $\mathbf{y}$  and  $\mathbf{z}$  are not. Assume we know that in this case, the best linear bounded  $\mathcal{MSE}$  predictor of  $\mathbf{y}$  and  $\mathbf{z}$  is  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$ , with prediction errors distributed as follows:

$$(20.1.17) \quad \begin{bmatrix} \mathbf{y} - \boldsymbol{\mu} \\ \mathbf{z} - \boldsymbol{\nu} \end{bmatrix} \sim \begin{bmatrix} \mathbf{o} \\ \mathbf{o} \end{bmatrix}, \sigma^2 \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix}.$$

This is the initial information. Here it is unnecessary to specify the unconditional distributions of  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$ , i.e.,  $\mathcal{E}[\boldsymbol{\mu}]$  and  $\mathcal{E}[\boldsymbol{\nu}]$  as well as the joint covariance matrix of  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  are not needed, even if they are known.

Then in a second step assume that an observation of  $\mathbf{y}$  becomes available, i.e., now  $\mathbf{y}$ ,  $\boldsymbol{\nu}$ , and  $\boldsymbol{\mu}$  are observed, but  $\mathbf{z}$  still isn't. Then the predictor

$$(20.1.18) \quad \mathbf{z}^* = \boldsymbol{\nu} + \Omega_{zy}\Omega_{yy}^-(\mathbf{y} - \boldsymbol{\mu})$$

is the best linear bounded  $\mathcal{MSE}$  predictor of  $\mathbf{z}$  based on  $\mathbf{y}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\nu}$ .

- a. Give special cases of this specification in which  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  are constant and  $\mathbf{y}$  and  $\mathbf{z}$  random, and one in which  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  and  $\mathbf{y}$  are random and  $\mathbf{z}$  is constant and one in which  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  are random and  $\mathbf{y}$  and  $\mathbf{z}$  are constant.

ANSWER. If  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  are constant, they are written  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$ . From this follows  $\boldsymbol{\mu} = \mathcal{E}[\boldsymbol{\mu}]$  and  $\boldsymbol{\nu} = \mathcal{E}[\boldsymbol{\nu}]$  and  $\sigma^2 \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix} = \mathcal{V} \begin{bmatrix} \mathbf{y} \\ r\mathbf{z} \end{bmatrix}$  and every linear predictor has bounded  $\mathcal{MSE}$ . Then the proof is as given earlier in this chapter. But an example in which  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  are not known constants but are observed random variables, and  $\mathbf{y}$  is also a random variable but  $\mathbf{z}$  is constant, is (21.0.2). Another example, in which  $\mathbf{y}$  and  $\mathbf{z}$  both are constants and  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  random, is constrained least squares (22.4.3).

- b. Prove equation 20.1.18.

ANSWER. In this proof we allow all four  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  and  $\mathbf{y}$  and  $\mathbf{z}$  to be random. A linear predictor based on  $\mathbf{y}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\nu}$  can be written as  $\tilde{\mathbf{z}} = \mathbf{B}\mathbf{y} + \mathbf{C}\boldsymbol{\mu} + \mathbf{D}\boldsymbol{\nu} + \mathbf{d}$ , therefore  $\tilde{\mathbf{z}} - \mathbf{z} = \mathbf{B}(\mathbf{y} - \boldsymbol{\mu}) + (\mathbf{C} + \mathbf{B})\boldsymbol{\mu} + (\mathbf{D} - \mathbf{I})\boldsymbol{\nu} - (\mathbf{z} - \boldsymbol{\nu}) + \mathbf{d}$ .  $\mathcal{E}[\tilde{\mathbf{z}} - \mathbf{z}] = \mathbf{o} + (\mathbf{C} + \mathbf{B})\mathcal{E}[\boldsymbol{\mu}] + (\mathbf{D} - \mathbf{I})\mathcal{E}[\boldsymbol{\nu}] - (\mathbf{z} - \boldsymbol{\nu}) + \mathbf{d}$ . Assuming that  $\mathcal{E}[\boldsymbol{\mu}]$  and  $\mathcal{E}[\boldsymbol{\nu}]$  can be anything, the requirement of bounded  $\mathcal{MSE}$  (or simply the requirement of unbiasedness, but this is not as elegant) gives  $\mathbf{C} = -\mathbf{B}$  and  $\mathbf{D} = \mathbf{I}$ , therefore  $\tilde{\mathbf{z}} = \boldsymbol{\nu} + \mathbf{B}(\mathbf{y} - \boldsymbol{\mu}) + \mathbf{d}$ , and the estimation error is  $\tilde{\mathbf{z}} - \mathbf{z} = \mathbf{B}(\mathbf{y} - \boldsymbol{\mu}) - (\mathbf{z} - \boldsymbol{\nu}) + \mathbf{d}$ . Now continue as in the proof of theorem 20.1.1. I must still carry out this proof much more carefully!

PROBLEM 269. 4 points According to (20.1.2), the prediction error  $\mathbf{z}^* - \mathbf{z}$  is uncorrelated with  $\mathbf{y}$ . If the distribution is such that the prediction error is independent of  $\mathbf{y}$  (as is the case if  $\mathbf{y}$  and  $\mathbf{z}$  are jointly normal), then  $\mathbf{z}^*$  as defined in (20.1.6) is the conditional mean  $\mathbf{z}^* = \mathcal{E}[\mathbf{z}|\mathbf{y}]$ , and its  $\mathcal{MSE}$ -matrix as defined in (20.1.7) is the conditional variance  $\mathcal{V}[\mathbf{z}|\mathbf{y}]$ .

ANSWER. From independence follows  $\mathcal{E}[\mathbf{z}^* - \mathbf{z}|\mathbf{y}] = \mathcal{E}[\mathbf{z}^* - \mathbf{z}]$ , and by the law of iterated expectations  $\mathcal{E}[\mathbf{z}^* - \mathbf{z}] = \mathbf{o}$ . Rewrite this as  $\mathcal{E}[\mathbf{z}|\mathbf{y}] = \mathcal{E}[\mathbf{z}^*|\mathbf{y}]$ . But since  $\mathbf{z}^*$  is a function of  $\mathbf{y}$ ,  $\mathcal{E}[\mathbf{z}^*|\mathbf{y}] = \mathbf{z}^*$ . Now the proof that the conditional dispersion matrix is the  $\mathcal{MSE}$  matrix:

$$(20.1.19) \quad \begin{aligned} \mathcal{V}[\mathbf{z}|\mathbf{y}] &= \mathcal{E}[(\mathbf{z} - \mathcal{E}[\mathbf{z}|\mathbf{y}])(\mathbf{z} - \mathcal{E}[\mathbf{z}|\mathbf{y}])^\top|\mathbf{y}] = \mathcal{E}[(\mathbf{z} - \mathbf{z}^*)(\mathbf{z} - \mathbf{z}^*)^\top|\mathbf{y}] \\ &= \mathcal{E}[(\mathbf{z} - \mathbf{z}^*)(\mathbf{z} - \mathbf{z}^*)^\top] = \mathcal{MSE}[\mathbf{z}^*; \mathbf{z}]. \end{aligned}$$

PROBLEM 270. Assume the expected values of  $\mathbf{x}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  are known, and the joint covariance matrix is known up to an unknown scalar factor  $\sigma^2 > 0$ .

$$(20.1.20) \quad \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\mu} \\ \boldsymbol{\nu} \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{\lambda} \\ \boldsymbol{\mu} \\ \boldsymbol{\nu} \end{bmatrix}, \sigma^2 \begin{bmatrix} \Omega_{xx} & \Omega_{xy} & \Omega_{xz} \\ \Omega_{xy}^\top & \Omega_{yy} & \Omega_{yz} \\ \Omega_{xz}^\top & \Omega_{yz}^\top & \Omega_{zz} \end{bmatrix}.$$

$\mathbf{x}$  is the original information,  $\mathbf{y}$  is additional information which becomes available, and  $\mathbf{z}$  is the variable which we want to predict on the basis of this information.

• a. 2 points Show that  $\mathbf{y}^* = \boldsymbol{\mu} + \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1}(\mathbf{x} - \boldsymbol{\lambda})$  is the best linear predictor of  $\mathbf{y}$  and  $\mathbf{z}^* = \boldsymbol{\nu} + \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1}(\mathbf{x} - \boldsymbol{\lambda})$  the best linear predictor of  $\mathbf{z}$  on the basis of the observation of  $\mathbf{x}$ , and that their joint MSE-matrix is

$$\mathcal{E} \left[ \begin{array}{c} \mathbf{y}^* - \mathbf{y} \\ \mathbf{z}^* - \mathbf{z} \end{array} \right] \left[ (\mathbf{y}^* - \mathbf{y})^\top \quad (\mathbf{z}^* - \mathbf{z})^\top \right] = \sigma^2 \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}} \\ \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}} & \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}} - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}} \end{bmatrix}$$

which can also be written

$$= \sigma^2 \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}|\mathbf{x}} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}|\mathbf{x}} \\ \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}|\mathbf{x}} & \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}|\mathbf{x}} \end{bmatrix}.$$

ANSWER. This part of the question is a simple application of the formulas derived earlier. For the MSE-matrix you first get

$$\sigma^2 \left( \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} \\ \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}^\top & \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}^\top \\ \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top \end{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}} \end{bmatrix} \right)$$

□

• b. 5 points Show that the best linear predictor of  $\mathbf{z}$  on the basis of the observations of  $\mathbf{x}$  and  $\mathbf{y}$  has the form

$$(20.1.21) \quad \mathbf{z}^{**} = \mathbf{z}^* + \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}|\mathbf{x}}^\top \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}|\mathbf{x}}^{-1}(\mathbf{y} - \mathbf{y}^*)$$

This is an important formula. All you need to compute  $\mathbf{z}^{**}$  is the best estimate  $\mathbf{z}^*$  before the new information  $\mathbf{y}$  became available, the best estimate  $\mathbf{y}^*$  of that new information itself, and the joint MSE matrix of the two. The original data  $\mathbf{x}$  and the covariance matrix (20.1.20) do not enter this formula.

ANSWER. Follows from

$$\mathbf{z}^{**} = \boldsymbol{\nu} + \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}^\top \end{bmatrix} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}} & \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \\ \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}^\top & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x} - \boldsymbol{\lambda} \\ \mathbf{y} - \boldsymbol{\mu} \end{bmatrix} =$$

Now apply (A.8.2):

$$\begin{aligned} &= \boldsymbol{\nu} + \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}^\top \end{bmatrix} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} + \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} & -\boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}} \\ -\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}} \end{bmatrix} \begin{bmatrix} \mathbf{x} - \boldsymbol{\lambda} \\ \mathbf{y} - \boldsymbol{\mu} \end{bmatrix} = \\ &= \boldsymbol{\nu} + \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}^\top \end{bmatrix} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1}(\mathbf{x} - \boldsymbol{\lambda}) + \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y}^* - \boldsymbol{\mu}) - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y} - \boldsymbol{\mu}) \\ -\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}(\mathbf{y}^* - \boldsymbol{\mu}) + \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y} - \boldsymbol{\mu}) \end{bmatrix} = \\ &= \boldsymbol{\nu} + \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}^\top \end{bmatrix} \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1}(\mathbf{x} - \boldsymbol{\lambda}) - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y} - \mathbf{y}^*) \\ +\boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y} - \mathbf{y}^*) \end{bmatrix} = \\ &= \boldsymbol{\nu} + \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1}(\mathbf{x} - \boldsymbol{\lambda}) - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y} - \mathbf{y}^*) + \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}^\top \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y} - \mathbf{y}^*) = \\ &= \mathbf{z}^* + (\boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}}^\top - \boldsymbol{\Omega}_{\mathbf{x}\mathbf{z}}^\top \boldsymbol{\Omega}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\Omega}_{\mathbf{x}\mathbf{y}}) \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{x}}(\mathbf{y} - \mathbf{y}^*) = \mathbf{z}^* + \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}|\mathbf{x}}^\top \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}|\mathbf{x}}^{-1}(\mathbf{y} - \mathbf{y}^*) \end{aligned}$$

PROBLEM 271. Assume  $\mathbf{x}$ ,  $\mathbf{y}$ , and  $\mathbf{z}$  have a joint probability distribution, and conditional expectation  $\mathcal{E}[\mathbf{z}|\mathbf{x}, \mathbf{y}] = \boldsymbol{\alpha}^* + \mathbf{A}^* \mathbf{x} + \mathbf{B}^* \mathbf{y}$  is linear in  $\mathbf{x}$  and  $\mathbf{y}$ .

• a. 1 point Show that  $\mathcal{E}[\mathbf{z}|\mathbf{x}] = \boldsymbol{\alpha}^* + \mathbf{A}^* \mathbf{x} + \mathbf{B}^* \mathcal{E}[\mathbf{y}|\mathbf{x}]$ . Hint: you may use law of iterated expectations in the following form:  $\mathcal{E}[\mathbf{z}|\mathbf{x}] = \mathcal{E}[\mathcal{E}[\mathbf{z}|\mathbf{x}, \mathbf{y}|\mathbf{x}]]$ .

ANSWER. With this hint it is trivial:  $\mathcal{E}[\mathbf{z}|\mathbf{x}] = \mathcal{E}[\boldsymbol{\alpha}^* + \mathbf{A}^* \mathbf{x} + \mathbf{B}^* \mathbf{y}|\mathbf{x}] = \boldsymbol{\alpha}^* + \mathbf{A}^* \mathbf{x} + \mathbf{B}^* \mathcal{E}[\mathbf{y}|\mathbf{x}]$ .

• b. 1 point The next three examples are from [CW99, pp. 264/5]: Assume  $\mathcal{E}[\mathbf{z}|\mathbf{x}, \mathbf{y}] = 1 + 2\mathbf{x} + 3\mathbf{y}$ ,  $\mathbf{x}$  and  $\mathbf{y}$  are independent, and  $\mathcal{E}[\mathbf{y}] = 2$ . Compute  $\mathcal{E}[\mathbf{z}|\mathbf{x}]$ .

ANSWER. According to the formula,  $\mathcal{E}[\mathbf{z}|\mathbf{x}] = 1 + 2\mathbf{x} + 3\mathcal{E}[\mathbf{y}|\mathbf{x}]$ , but since  $\mathbf{x}$  and  $\mathbf{y}$  are independent  $\mathcal{E}[\mathbf{y}|\mathbf{x}] = \mathcal{E}[\mathbf{y}] = 2$ ; therefore  $\mathcal{E}[\mathbf{z}|\mathbf{x}] = 7 + 2\mathbf{x}$ . I.e., the slope is the same, but the intercept changes.

• c. 1 point Assume again  $\mathcal{E}[\mathbf{z}|\mathbf{x}, \mathbf{y}] = 1 + 2\mathbf{x} + 3\mathbf{y}$ , but this time  $\mathbf{x}$  and  $\mathbf{y}$  are independent but  $\mathcal{E}[\mathbf{y}|\mathbf{x}] = 2 - \mathbf{x}$ . Compute  $\mathcal{E}[\mathbf{z}|\mathbf{x}]$ .

ANSWER.  $\mathcal{E}[\mathbf{z}|\mathbf{x}] = 1 + 2\mathbf{x} + 3(2 - \mathbf{x}) = 7 - \mathbf{x}$ . In this situation, both slope and intercept change, but it is still a linear relationship.

• d. 1 point Again  $\mathcal{E}[\mathbf{z}|\mathbf{x}, \mathbf{y}] = 1 + 2\mathbf{x} + 3\mathbf{y}$ , and this time the relationship between  $\mathbf{x}$  and  $\mathbf{y}$  is nonlinear:  $\mathcal{E}[\mathbf{y}|\mathbf{x}] = 2 - e^{\mathbf{x}}$ . Compute  $\mathcal{E}[\mathbf{z}|\mathbf{x}]$ .

ANSWER.  $\mathcal{E}[\mathbf{z}|\mathbf{x}] = 1 + 2\mathbf{x} + 3(2 - e^{\mathbf{x}}) = 7 + 2\mathbf{x} - 3e^{\mathbf{x}}$ . This time the marginal relationship between  $\mathbf{x}$  and  $\mathbf{y}$  is no longer linear. This is so despite the fact that, if all the variables are included, i.e., if both  $\mathbf{x}$  and  $\mathbf{y}$  are included, then the relationship is linear.

• e. 1 point Assume  $\mathcal{E}[f(\mathbf{z})|\mathbf{x}, \mathbf{y}] = 1 + 2\mathbf{x} + 3\mathbf{y}$ , where  $f$  is a nonlinear function and  $\mathcal{E}[\mathbf{y}|\mathbf{x}] = 2 - \mathbf{x}$ . Compute  $\mathcal{E}[f(\mathbf{z})|\mathbf{x}]$ .

ANSWER.  $\mathcal{E}[f(\mathbf{z})|\mathbf{x}] = 1 + 2\mathbf{x} + 3(2 - \mathbf{x}) = 7 - \mathbf{x}$ . If one plots  $\mathbf{z}$  against  $\mathbf{x}$  and  $\mathbf{z}$ , then the plots should be similar, though not identical, since the same transformation  $f$  will straighten them out. This is why the plots in the top row or right column of [CW99, p. 435] are so similar.

Connection between prediction and inverse prediction: If  $\mathbf{y}$  is observed and  $\mathbf{z}$  is to be predicted, the BLUP is  $\mathbf{z}^* - \boldsymbol{\nu} = \mathbf{B}^*(\mathbf{y} - \boldsymbol{\mu})$  where  $\mathbf{B}^* = \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}}^{-1}$ . If  $\mathbf{z}$  is observed and  $\mathbf{y}$  is to be predicted, then the BLUP is  $\mathbf{y}^* - \boldsymbol{\mu} = \mathbf{C}^*(\mathbf{z} - \boldsymbol{\nu})$  where  $\mathbf{C}^* = \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}}^{-1}$ .  $\mathbf{B}^*$  and  $\mathbf{C}^*$  are connected by the formula

$$(20.1.22) \quad \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} \mathbf{B}^{*\top} = \mathbf{C}^* \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}}.$$

This relationship can be used for graphical regression methods [Coo98, pp. 187/188]. If  $\mathbf{z}$  is a scalar, it is much easier to determine the elements of  $\mathbf{C}^*$  than those of  $\mathbf{B}^*$ .  $\mathbf{C}^*$  consists of the regression slopes in the scatter plot of each of the observed variables against  $\mathbf{z}$ . They can be read off easily from a scatterplot matrix. This works not only if the distribution is Normal, but also with arbitrary distributions as long as all conditional expectations between the explanatory variables are linear.

PROBLEM 272. In order to make relationship (20.1.22) more intuitive, assume  $x$  and  $\varepsilon$  are Normally distributed and independent of each other, and  $E[\varepsilon] = 0$ . Define  $y = \alpha + \beta x + \varepsilon$ .

• a. Show that  $\alpha + \beta x$  is the best linear predictor of  $y$  based on the observation of  $x$ .

ANSWER. Follows from the fact that the predictor is unbiased and the prediction error is uncorrelated with  $x$ .  $\square$

• b. Express  $\beta$  in terms of the variances and covariances of  $x$  and  $y$ .

ANSWER.  $\text{cov}[x, y] = \beta \text{var}[x]$ , therefore  $\beta = \frac{\text{cov}[x, y]}{\text{var}[x]}$   $\square$

• c. Since  $x$  and  $y$  are jointly normal, they can also be written  $x = \gamma + \delta y + \omega$  where  $\omega$  is independent of  $y$ . Express  $\delta$  in terms of the variances and covariances of  $x$  and  $y$ , and show that  $\text{var}[y]\beta = \gamma \text{var}[x]$ .

ANSWER.  $\delta = \frac{\text{cov}[x, y]}{\text{var}[y]}$ .  $\square$

• d. Now let us extend the model a little: assume  $x_1, x_2$ , and  $\varepsilon$  are Normally distributed and independent of each other, and  $E[\varepsilon] = 0$ . Define  $y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$ . Again express  $\beta_1$  and  $\beta_2$  in terms of variances and covariances of  $x_1, x_2$ , and  $y$ .

ANSWER. Since  $x_1$  and  $x_2$  are independent, one gets the same formulas as in the univariate case: from  $\text{cov}[x_1, y] = \beta_1 \text{var}[x_1]$  and  $\text{cov}[x_2, y] = \beta_2 \text{var}[x_2]$  follows  $\beta_1 = \frac{\text{cov}[x_1, y]}{\text{var}[x_1]}$  and  $\beta_2 = \frac{\text{cov}[x_2, y]}{\text{var}[x_2]}$ .  $\square$

• e. Since  $x_1$  and  $y$  are jointly normal, they can also be written  $x_1 = \gamma_1 + \delta_1 y + \omega_1$ , where  $\omega_1$  is independent of  $y$ . Likewise,  $x_2 = \gamma_2 + \delta_2 y + \omega_2$ , where  $\omega_2$  is independent of  $y$ . Express  $\delta_1$  and  $\delta_2$  in terms of the variances and covariances of  $x_1, x_2$ , and  $y$ , and show that

$$(20.1.23) \quad \begin{bmatrix} \delta_1 \\ \delta_2 \end{bmatrix} \text{var}[y] = \begin{bmatrix} \text{var}[x_1] & 0 \\ 0 & \text{var}[x_2] \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

This is (20.1.22) in the present situation.

ANSWER.  $\delta_1 = \frac{\text{cov}[x_1, y]}{\text{var}[y]}$  and  $\delta_2 = \frac{\text{cov}[x_2, y]}{\text{var}[y]}$ .  $\square$

### 20.2. The Associated Least Squares Problem

For every estimation problem there is an associated “least squares” problem. In the present situation,  $z^*$  is that value which, together with the given observation  $y$ , “blends best” into the population defined by  $\mu, \nu$  and the dispersion matrix  $\Omega$ , in the following sense: Given the observed value  $y$ , the vector  $z^* = \nu + \Omega_{zy} \Omega_{yy}^{-1} (y - \mu)$

is that value  $z$  for which  $\begin{bmatrix} y \\ z \end{bmatrix}$  has smallest Mahalanobis distance from the population

defined by the mean vector  $\begin{bmatrix} \mu \\ \nu \end{bmatrix}$  and the covariance matrix  $\sigma^2 \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix}$ .

In the case of singular  $\Omega_{zz}$ , it is only necessary to minimize among those which have finite distance from the population, i.e., which can be written in the form  $z = \nu + \Omega_{zz} q$  for some  $q$ . We will also write  $r = \text{rank} \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix}$ . Therefore, it solves the following “least squares problem:”

$$(20.2.1) \quad z = z^* \quad \min. \frac{1}{r\sigma^2} \begin{bmatrix} y - \mu \\ z - \nu \end{bmatrix}^\top \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix}^{-1} \begin{bmatrix} y - \mu \\ z - \nu \end{bmatrix} \quad \text{s. t. } z = \nu + \Omega_{zz} q \text{ for some } q$$

To prove this, use (A.8.2) to invert the dispersion matrix:

$$(20.2.2) \quad \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{zy} & \Omega_{zz} \end{bmatrix}^{-1} = \begin{bmatrix} \Omega_{yy}^{-1} + \Omega_{yy}^{-1} \Omega_{yz} \Omega_{zz}^{-1} \Omega_{zy} \Omega_{yy}^{-1} & -\Omega_{yy}^{-1} \Omega_{yz} \Omega_{zz}^{-1} \\ -\Omega_{zz}^{-1} \Omega_{zy} \Omega_{yy}^{-1} & \Omega_{zz}^{-1} \end{bmatrix}.$$

If one plugs  $z = z^*$  into this objective function, one obtains a very simple expression:

$$(20.2.3) \quad (y - \mu)^\top \begin{bmatrix} I & \Omega_{yy}^{-1} \Omega_{yz} \\ & -\Omega_{zz}^{-1} \Omega_{zy} \Omega_{yy}^{-1} \end{bmatrix} \begin{bmatrix} \Omega_{yy}^{-1} + \Omega_{yy}^{-1} \Omega_{yz} \Omega_{zz}^{-1} \Omega_{zy} \Omega_{yy}^{-1} & -\Omega_{yy}^{-1} \Omega_{yz} \Omega_{zz}^{-1} \\ & \Omega_{zz}^{-1} \end{bmatrix} \begin{bmatrix} I \\ \Omega_{zy} \Omega_{yy}^{-1} \end{bmatrix} (y - \mu)$$

$$(20.2.4) \quad = (y - \mu)^\top \Omega_{yy}^{-1} (y - \mu).$$

Now take any  $z$  of the form  $z = \nu + \Omega_{zz} q$  for some  $q$  and write it in the form  $z = z^* + \Omega_{zz} d$ , i.e.,

$$\begin{bmatrix} y - \mu \\ z - \nu \end{bmatrix} = \begin{bmatrix} y - \mu \\ z^* - \nu \end{bmatrix} + \begin{bmatrix} o \\ \Omega_{zz} d \end{bmatrix}.$$

Then the cross product terms in the objective function disappear:

$$(20.2.5) \quad \begin{bmatrix} o^\top & d^\top \Omega_{zz} \end{bmatrix} \begin{bmatrix} \Omega_{yy}^{-1} + \Omega_{yy}^{-1} \Omega_{yz} \Omega_{zz}^{-1} \Omega_{zy} \Omega_{yy}^{-1} & -\Omega_{yy}^{-1} \Omega_{yz} \Omega_{zz}^{-1} \\ -\Omega_{zz}^{-1} \Omega_{zy} \Omega_{yy}^{-1} & \Omega_{zz}^{-1} \end{bmatrix} \begin{bmatrix} I \\ \Omega_{zy} \Omega_{yy}^{-1} \end{bmatrix} (y - \mu) + \begin{bmatrix} o^\top & d^\top \Omega_{zz} \end{bmatrix} \begin{bmatrix} \Omega_{yy}^{-1} \\ o \end{bmatrix} (y - \mu) = \dots$$

Therefore this gives a larger value of the objective function.

PROBLEM 273. Use problem 379 for an alternative proof of this.

From (20.2.1) follows that  $z^*$  is the mode of the normal density function, and since the mode is the mean, this is an alternative proof, in the case of nonsingular covariance matrix, when the density exists, that  $z^*$  is the normal conditional mean.

**20.3. Prediction of Future Observations in the Regression Model**

For a moment let us go back to the model  $\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\varepsilon}$  with spherically distributed disturbances  $\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2\mathbf{I})$ . This time, our goal is not to estimate  $\beta$ , but the situation is the following: For a new set of observations of the explanatory variables  $\mathbf{X}_0$  the values of the dependent variable  $\mathbf{y}_0 = \mathbf{X}_0\beta + \boldsymbol{\varepsilon}_0$  have not yet been observed and we want to predict them. The obvious predictor is  $\mathbf{y}_0^* = \mathbf{X}_0\hat{\beta} = \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}$ .

Since

$$(20.3.1) \quad \begin{aligned} \mathbf{y}_0^* - \mathbf{y}_0 &= \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y} - \mathbf{y}_0 = \\ &= \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}\beta + \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\varepsilon} - \mathbf{X}_0\beta - \boldsymbol{\varepsilon}_0 = \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_0 \end{aligned}$$

one sees that  $E[\mathbf{y}_0^* - \mathbf{y}_0] = \mathbf{o}$ , i.e., it is an unbiased predictor. And since  $\boldsymbol{\varepsilon}$  and  $\boldsymbol{\varepsilon}_0$  are uncorrelated, one obtains

$$(20.3.2) \quad \mathcal{MSE}[\mathbf{y}_0^*; \mathbf{y}_0] = \mathcal{V}[\mathbf{y}_0^* - \mathbf{y}_0] = \mathcal{V}[\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\varepsilon}] + \mathcal{V}[\boldsymbol{\varepsilon}_0]$$

$$(20.3.3) \quad = \sigma^2(\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}_0^\top + \mathbf{I}).$$

Problem 274 shows that this is the Best Linear Unbiased Predictor (BLUP) of  $\mathbf{y}_0$  on the basis of  $\mathbf{y}$ .

PROBLEM 274. *The prediction problem in the Ordinary Least Squares model can be formulated as follows:*

$$(20.3.4) \quad \begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{X}_0 \end{bmatrix} \beta + \begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix} \quad \mathcal{E}\left[\begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix}\right] = \begin{bmatrix} \mathbf{o} \\ \mathbf{o} \end{bmatrix} \quad \mathcal{V}\left[\begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix}\right] = \sigma^2 \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix}.$$

$\mathbf{X}$  and  $\mathbf{X}_0$  are known,  $\mathbf{y}$  is observed,  $\mathbf{y}_0$  is not observed.

• a. 4 points Show that  $\mathbf{y}_0^* = \mathbf{X}_0\hat{\beta}$  is the Best Linear Unbiased Predictor (BLUP) of  $\mathbf{y}_0$  on the basis of  $\mathbf{y}$ , where  $\hat{\beta}$  is the OLS estimate in the model  $\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\varepsilon}$ .

ANSWER. Take any other predictor  $\tilde{\mathbf{y}}_0 = \tilde{\mathbf{B}}\mathbf{y}$  and write  $\tilde{\mathbf{B}} = \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D}$ . Unbiasedness means  $\mathcal{E}[\tilde{\mathbf{y}}_0 - \mathbf{y}_0] = \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}\beta + \mathbf{D}\mathbf{X}\beta - \mathbf{X}_0\beta = \mathbf{o}$ , from which follows  $\mathbf{D}\mathbf{X} = \mathbf{O}$ . Because of unbiasedness we know  $\mathcal{MSE}[\tilde{\mathbf{y}}_0; \mathbf{y}_0] = \mathcal{V}[\tilde{\mathbf{y}}_0 - \mathbf{y}_0]$ . Since the prediction error can be written  $\tilde{\mathbf{y}}_0 - \mathbf{y}_0 = [\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D} \quad -\mathbf{I}] \begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix}$ , one obtains

$$\begin{aligned} \mathcal{V}[\tilde{\mathbf{y}}_0 - \mathbf{y}_0] &= [\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D} \quad -\mathbf{I}] \mathcal{V}\left[\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix}\right] \begin{bmatrix} \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D}^\top \\ -\mathbf{I} \end{bmatrix} \\ &= \sigma^2 [\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D} \quad -\mathbf{I}] \begin{bmatrix} \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D}^\top \\ -\mathbf{I} \end{bmatrix} \\ &= \sigma^2 (\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D})(\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{D})^\top + \sigma^2\mathbf{I} \\ &= \sigma^2 (\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}_0^\top + \mathbf{D}\mathbf{D}^\top + \mathbf{I}). \end{aligned}$$

This is smallest for  $\mathbf{D} = \mathbf{O}$ . □

• b. 2 points From our formulation of the Gauss-Markov theorem in Theorem 18.1.1 it is obvious that the same  $\mathbf{y}_0^* = \mathbf{X}_0\hat{\beta}$  is also the Best Linear Unbiased Predictor of  $\mathbf{X}_0\beta$ , which is the expected value of  $\mathbf{y}_0$ . You are not required to prove this here, but you are asked to compute  $\mathcal{MSE}[\mathbf{X}_0\hat{\beta}; \mathbf{X}_0\beta]$  and compare it with  $\mathcal{MSE}[\mathbf{y}_0^*; \mathbf{y}_0]$ . Can you explain the difference?

ANSWER. Estimation error and  $\mathcal{MSE}$  are

$$\begin{aligned} \mathbf{X}_0\hat{\beta} - \mathbf{X}_0\beta &= \mathbf{X}_0(\hat{\beta} - \beta) = \mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\varepsilon} \quad \text{due to (??)} \\ \mathcal{MSE}[\mathbf{X}_0\hat{\beta}; \mathbf{X}_0\beta] &= \mathcal{V}[\mathbf{X}_0\hat{\beta} - \mathbf{X}_0\beta] = \mathcal{V}[\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\varepsilon}] = \sigma^2\mathbf{X}_0(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}_0^\top. \end{aligned}$$

It differs from the prediction  $\mathcal{MSE}$  matrix by  $\sigma^2\mathbf{I}$ , which is the uncertainty about the value of new disturbance  $\boldsymbol{\varepsilon}_0$  about which the data have no information.

[Gre97, p. 369] has an enlightening formula showing how the prediction interval increase if one goes away from the center of the data.

Now let us look at the prediction problem in the Generalized Least Squares model

$$(20.3.5) \quad \begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{X}_0 \end{bmatrix} \beta + \begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix} \quad \mathcal{E}\left[\begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix}\right] = \begin{bmatrix} \mathbf{o} \\ \mathbf{o} \end{bmatrix} \quad \mathcal{V}\left[\begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix}\right] = \sigma^2 \begin{bmatrix} \boldsymbol{\Psi} & \mathbf{C} \\ \mathbf{C}^\top & \boldsymbol{\Psi}_0 \end{bmatrix}$$

$\mathbf{X}$  and  $\mathbf{X}_0$  are known,  $\mathbf{y}$  is observed,  $\mathbf{y}_0$  is not observed, and we assume  $\boldsymbol{\Psi}$  is positive definite. If  $\mathbf{C} = \mathbf{O}$ , the BLUP of  $\mathbf{y}_0$  is  $\mathbf{X}_0\hat{\beta}$ , where  $\hat{\beta}$  is the BLUE in the model  $\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\varepsilon}$ . In other words, all new disturbances are simply predicted by zero. If past and future disturbances are correlated, this predictor is no longer optimal.

In [JHG<sup>+</sup>88, pp. 343–346] it is proved that the best linear unbiased predictor of  $\mathbf{y}_0$  is

$$(20.3.6) \quad \mathbf{y}_0^* = \mathbf{X}_0\hat{\beta} + \mathbf{C}^\top\boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta}).$$

where  $\hat{\beta}$  is the generalized least squares estimator of  $\beta$ , and that its  $\mathcal{MSE}$ -matrix  $\mathcal{MSE}[\mathbf{y}_0^*; \mathbf{y}_0]$  is

$$(20.3.7) \quad \sigma^2 \left( \boldsymbol{\Psi}_0 - \mathbf{C}^\top\boldsymbol{\Psi}^{-1}\mathbf{C} + (\mathbf{X}_0 - \mathbf{C}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})(\mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top\boldsymbol{\Psi}^{-1}\mathbf{C}) \right)$$

PROBLEM 275. *Derive the formula for the  $\mathcal{MSE}$  matrix from the formula for the predictor, and compute the joint  $\mathcal{MSE}$  matrix for the predicted values and parameter vector.*

ANSWER. The prediction error is, using (19.0.7),

$$(20.3.8) \quad \mathbf{y}_0^* - \mathbf{y}_0 = \mathbf{X}_0 \hat{\boldsymbol{\beta}} - \mathbf{X}_0 \boldsymbol{\beta} + \mathbf{X}_0 \boldsymbol{\beta} - \mathbf{y}_0 + \mathbf{C}^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} + \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

$$(20.3.9) \quad = \mathbf{X}_0(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) - \boldsymbol{\varepsilon}_0 + \mathbf{C}^\top \boldsymbol{\Psi}^{-1}(\boldsymbol{\varepsilon} - \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}))$$

$$(20.3.10) \quad = \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \boldsymbol{\varepsilon} + (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) - \boldsymbol{\varepsilon}_0$$

$$(20.3.11) \quad = \left[ \mathbf{C}^\top \boldsymbol{\Psi}^{-1} + (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \quad -\mathbf{I} \right] \begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix}$$

The  $MSE$ -matrix is therefore  $\mathcal{E}[(\mathbf{y}_0^* - \mathbf{y}_0)(\mathbf{y}_0^* - \mathbf{y}_0)^\top] =$

$$(20.3.12) \quad = \sigma^2 \left[ \mathbf{C}^\top \boldsymbol{\Psi}^{-1} + (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \quad -\mathbf{I} \right] \\ \begin{bmatrix} \boldsymbol{\Psi} & \mathbf{C} \\ \mathbf{C}^\top & \boldsymbol{\Psi}_0 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Psi}^{-1} \mathbf{C} + \boldsymbol{\Psi}^{-1} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C}) \\ -\mathbf{I} \end{bmatrix}$$

and the joint  $MSE$  matrix with the sampling error of the parameter vector  $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}$  is

$$(20.3.13) \quad \sigma^2 \begin{bmatrix} \mathbf{C}^\top \boldsymbol{\Psi}^{-1} + (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} & -\mathbf{I} \\ (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} & \mathbf{O} \end{bmatrix} \\ \begin{bmatrix} \boldsymbol{\Psi} & \mathbf{C} \\ \mathbf{C}^\top & \boldsymbol{\Psi}_0 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Psi}^{-1} \mathbf{C} + \boldsymbol{\Psi}^{-1} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C}) & \boldsymbol{\Psi}^{-1} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \\ -\mathbf{I} & \mathbf{O} \end{bmatrix} =$$

$$(20.3.14) \quad = \sigma^2 \begin{bmatrix} \mathbf{C}^\top \boldsymbol{\Psi}^{-1} + (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} & -\mathbf{I} \\ (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} & \mathbf{O} \end{bmatrix} \\ \begin{bmatrix} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C}) & \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \\ \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{C} + \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C}) - \boldsymbol{\Psi}_0 & \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \end{bmatrix}$$

If one multiplies this out, one gets

$$(20.3.15) \quad \begin{bmatrix} \boldsymbol{\Psi}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{C} + (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C}) & (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \\ (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C}) & (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \end{bmatrix}$$

The upper left diagonal element is as claimed in (20.3.7).  $\square$

The strategy of the proof given in ITPE is similar to the strategy used to obtain the GLS results, namely, to transform the data in such a way that the disturbances are well behaved. Both data vectors  $\mathbf{y}$  and  $\mathbf{y}_0$  will be transformed, but this transformation must have the following additional property: the transformed  $\mathbf{y}$  must be a function of  $\mathbf{y}$  alone, not of  $\mathbf{y}_0$ . Once such a transformation is found, it is easy to predict the transformed  $\mathbf{y}_0$  on the basis of the transformed  $\mathbf{y}$ , and from this one also obtains a prediction of  $\mathbf{y}_0$  on the basis of  $\mathbf{y}$ .

Here is some heuristics in order to understand formula (20.3.6). Assume for a moment that  $\boldsymbol{\beta}$  was known. Then you can apply theorem ?? to the model

$$(20.3.16) \quad \begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix} \sim \begin{bmatrix} \mathbf{X}\boldsymbol{\beta} \\ \mathbf{X}_0\boldsymbol{\beta} \end{bmatrix}, \sigma^2 \begin{bmatrix} \boldsymbol{\Psi} & \mathbf{C} \\ \mathbf{C}^\top & \boldsymbol{\Psi}_0 \end{bmatrix}$$

to get  $\mathbf{y}_0^* = \mathbf{X}_0 \boldsymbol{\beta} + \mathbf{C}^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  as best linear predictor of  $\mathbf{y}_0$  on the basis  $\mathbf{y}$ . According to theorem ??, its  $MSE$  matrix is  $\sigma^2(\boldsymbol{\Psi}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{C})$ . Since  $\boldsymbol{\beta}$  not known, replace it by  $\hat{\boldsymbol{\beta}}$ , which gives exactly (20.3.6). This adds  $MSE[\mathbf{X}_0 \hat{\boldsymbol{\beta}} + \mathbf{C}^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}); \mathbf{X}_0 \boldsymbol{\beta} + \mathbf{C}^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})]$  to the  $MSE$ -matrix, which gives (20.3.12).

PROBLEM 276. Show that

$$(20.3.17) \quad MSE[\mathbf{X}_0 \hat{\boldsymbol{\beta}} + \mathbf{C}^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}); \mathbf{X}_0 \boldsymbol{\beta} + \mathbf{C}^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})] = \\ = \sigma^2(\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C})$$

ANSWER. What is predicted is a random variable, therefore the  $MSE$  matrix is the covariance matrix of the prediction error. The prediction error is  $(\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ , its covariance matrix is therefore  $\sigma^2(\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1}(\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C})$ .

PROBLEM 277. In the following we work with partitioned matrices. Given a model

$$(20.3.18) \quad \begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{X}_0 \end{bmatrix} \boldsymbol{\beta} + \begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix} \quad \mathcal{E} \begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{o} \\ \mathbf{o} \end{bmatrix} \quad \mathcal{V} \begin{bmatrix} \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon}_0 \end{bmatrix} = \sigma^2 \begin{bmatrix} \boldsymbol{\Psi} & \mathbf{C} \\ \mathbf{C}^\top & \boldsymbol{\Psi}_0 \end{bmatrix}$$

$\mathbf{X}$  has full rank.  $\mathbf{y}$  is observed,  $\mathbf{y}_0$  is not observed.  $\mathbf{C}$  is not the null matrix.

• a. Someone predicts  $\mathbf{y}_0$  by  $\mathbf{y}_0^* = \mathbf{X}_0 \hat{\boldsymbol{\beta}}$ , where  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{y}$  is the BLUE of  $\boldsymbol{\beta}$ . Is this predictor unbiased?

ANSWER. Yes, since  $\mathcal{E}[\mathbf{y}_0] = \mathbf{X}_0 \boldsymbol{\beta}$ , and  $\mathcal{E}[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ .

• b. Compute the  $MSE$  matrix  $MSE[\mathbf{X}_0 \hat{\boldsymbol{\beta}}; \mathbf{y}_0]$  of this predictor. Hint: For a partitioned matrix  $\mathbf{B}$ , the difference  $\mathbf{B}\mathbf{y} - \mathbf{y}_0$  can be written in the form  $[\mathbf{B} \quad -\mathbf{I}] \begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix}$ . Hint: Use (20.3.18).

For an unbiased predictor (or estimator), the  $MSE$  matrix is the covariance matrix of the prediction (or estimation) error.

ANSWER.

$$(20.3.19) \quad \mathcal{E}[(\mathbf{B}\mathbf{y} - \mathbf{y}_0)(\mathbf{B}\mathbf{y} - \mathbf{y}_0)^\top] = \mathcal{V}[\mathbf{B}\mathbf{y} - \mathbf{y}_0]$$

$$(20.3.20) \quad = \mathcal{V} \left[ [\mathbf{B} \quad -\mathbf{I}] \begin{bmatrix} \mathbf{y} \\ \mathbf{y}_0 \end{bmatrix} \right]$$

$$(20.3.21) \quad = \sigma^2 [\mathbf{B} \quad -\mathbf{I}] \begin{bmatrix} \boldsymbol{\Psi} & \mathbf{C} \\ \mathbf{C}^\top & \boldsymbol{\Psi}_0 \end{bmatrix} \begin{bmatrix} \mathbf{B}^\top \\ -\mathbf{I} \end{bmatrix}$$

$$(20.3.22) \quad = \sigma^2 (\mathbf{B}\boldsymbol{\Psi}\mathbf{B}^\top - \mathbf{C}^\top \mathbf{B}^\top - \mathbf{C}\mathbf{B} + \boldsymbol{\Psi}_0).$$

Now one must use  $\mathbf{B} = \mathbf{X}_0(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1}$ . One ends up with

$$(20.3.23) \quad MSE[\mathbf{X}_0 \hat{\boldsymbol{\beta}}; \mathbf{y}_0] = \sigma^2 \left( \mathbf{X}_0(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}_0^\top - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}_0^\top - \mathbf{X}_0(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C} + \boldsymbol{\Psi}_0 \right)$$



• c. Compare its  $\mathcal{MSE}$ -matrix with formula (20.3.7). Is the difference nonnegative definite?

ANSWER. To compare it with the minimum  $\mathcal{MSE}$  matrix, it can also be written as (20.3.24)

$$\mathcal{MSE}[\mathbf{X}_0 \hat{\boldsymbol{\beta}}; \mathbf{y}_0] = \sigma^2 \left( \boldsymbol{\Psi}_0 + (\mathbf{X}_0 - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} (\mathbf{X}_0^\top - \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C}) - \mathbf{C}^\top \boldsymbol{\Psi}^{-1} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{C} \right)$$

i.e., it exceeds the minimum  $\mathcal{MSE}$  matrix by  $\mathbf{C}^\top (\boldsymbol{\Psi}^{-1} - \boldsymbol{\Psi}^{-1} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1}) \mathbf{C}$ . This is nnd because the matrix in parentheses is  $\mathbf{M} = \mathbf{M} \boldsymbol{\Psi} \mathbf{M}$ , refer here to Problem 265.  $\square$

## CHAPTER 21

## Updating of Estimates When More Observations become Available

The theory of the linear model often deals with pairs of models which are nested in each other, one model either having more data or more stringent parameter restrictions than the other. We will discuss such nested models in three forms: in the remainder of the present chapter 21 we will see how estimates must be updated when more observations become available, in chapter 22 how the imposition of a linear constraint affects the parameter estimates, and in chapter 23 what happens if one adds more regressors.

Assume you have already computed the BLUE  $\hat{\beta}$  on the basis of the observations  $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ , and afterwards additional data  $\mathbf{y}_0 = \mathbf{X}_0\beta + \varepsilon_0$  become available. Then  $\hat{\beta}$  can be updated using the following principles:

Before the new observations became available, the information given in the original dataset not only allowed to estimate  $\beta$  by  $\hat{\beta}$ , but also yielded a prediction  $\mathbf{y}_0^* = \mathbf{X}_0\hat{\beta}$  of the additional data. The estimation error  $\hat{\beta} - \beta$  and the prediction error  $\mathbf{y}_0^* - \mathbf{y}_0$  are unobserved, but we know their expected values (the zero vectors), and we also know their joint covariance matrix up to the unknown factor  $\sigma^2$ . After the additional data have become available, we can compute the actual value of the prediction error  $\mathbf{y}_0^* - \mathbf{y}_0$ . This allows us to also get a better idea of the actual value of the estimation error, and therefore we can get a better estimator of  $\beta$ . The following steps are involved:

- (1) Make the best prediction  $\mathbf{y}_0^*$  of the new data  $\mathbf{y}_0$  based on  $\hat{\beta}$ .
- (2) Compute the joint covariance matrix of the prediction error  $\mathbf{y}_0^* - \mathbf{y}_0$  of the new data by the old (which is observed) and the sampling error in the old regression  $\hat{\beta} - \beta$  (which is unobserved).
- (3) Use the formula for best linear prediction (??) to get a predictor  $\mathbf{z}^*$  of  $\hat{\beta} - \beta$ .
- (4) Then  $\hat{\beta} = \hat{\beta} - \mathbf{z}^*$  is the BLUE of  $\beta$  based on the joint observations  $\mathbf{y}$  and  $\mathbf{y}_0$ .

(5) The sum of squared errors of the updated model minus that of the base model is the standardized prediction error  $\text{SSE}^* - \text{SSE} = (\mathbf{y}_0^* - \mathbf{y}_0)^\top \mathbf{\Omega}^{-1} (\mathbf{y}_0^* - \mathbf{y}_0)$  where  $\text{SSE}^* = (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) \mathcal{V}[\mathbf{y}_0^* - \mathbf{y}_0] = \sigma^2 \mathbf{\Omega}$ .

In the case of *one* additional observation and spherical covariance matrix, the procedure yields the following formulas:

PROBLEM 278. Assume  $\hat{\beta}$  is the BLUE on the basis of the observation  $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ , and a new observation  $\mathbf{y}_0 = \mathbf{x}_0^\top \beta + \varepsilon_0$  becomes available. Show that the updated estimator has the form

$$(21.0.25) \quad \hat{\beta} = \hat{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 \frac{\mathbf{y}_0 - \mathbf{x}_0^\top \hat{\beta}}{1 + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0}.$$

ANSWER. Set it up as follows:

$$(21.0.26) \quad \begin{bmatrix} \mathbf{y}_0 - \mathbf{x}_0^\top \hat{\beta} \\ \hat{\beta} - \beta \end{bmatrix} \sim \begin{bmatrix} 0 \\ \sigma \end{bmatrix}, \sigma^2 \begin{bmatrix} \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 + 1 & \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \\ (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 & (\mathbf{X}^\top \mathbf{X})^{-1} \end{bmatrix}$$

and use (20.1.18). By the way, if the covariance matrix is not spherical but is  $\begin{bmatrix} \Psi & c \\ c^\top & \psi_0 \end{bmatrix}$  we from (20.3.6)

$$(21.0.27) \quad \mathbf{y}_0^* = \mathbf{x}_0^\top \hat{\beta} + \mathbf{c}^\top \Psi^{-1} (\mathbf{y} - \mathbf{X}\hat{\beta})$$

and from (20.3.15)

$$(21.0.28) \quad \begin{bmatrix} \mathbf{y}_0 - \mathbf{y}_0^* \\ \hat{\beta} - \beta \end{bmatrix} \sim \begin{bmatrix} 0 \\ \sigma \end{bmatrix}, \sigma^2 \begin{bmatrix} \psi_0 - \mathbf{c}^\top \Psi^{-1} \mathbf{c} + (\mathbf{x}_0^\top - \mathbf{c}^\top \Psi^{-1} \mathbf{X}) (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} (\mathbf{x}_0 - \mathbf{X}^\top \Psi^{-1} \mathbf{c}) & (\mathbf{x}_0^\top - \mathbf{c}^\top \Psi^{-1} \mathbf{X}) (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} \\ (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} (\mathbf{x}_0 - \mathbf{X}^\top \Psi^{-1} \mathbf{c}) & (\mathbf{X}^\top \Psi^{-1} \mathbf{X})^{-1} \end{bmatrix}$$

• a. Show that the residual  $\hat{\varepsilon}_0$  from the full regression is the following nonrandom multiple of the “predictive” residual  $\mathbf{y}_0 - \mathbf{x}_0^\top \hat{\beta}$ :

$$(21.0.29) \quad \hat{\varepsilon}_0 = \mathbf{y}_0 - \mathbf{x}_0^\top \hat{\beta} = \frac{1}{1 + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0} (\mathbf{y}_0 - \mathbf{x}_0^\top \hat{\beta})$$

Interestingly, this is the predictive residual divided by its relative variance (to standardize it one would have to divide it by its relative standard deviation). Compute this with (24.2.9).

ANSWER. (21.0.29) can either be derived from (21.0.25), or from the following alternative application of the updating principle: All the information which the old observations have for estimate of  $\mathbf{x}_0^\top \beta$  is contained in  $\hat{\mathbf{y}}_0 = \mathbf{x}_0^\top \hat{\beta}$ . The information which the updated regression, with

includes the additional observation, has about  $\mathbf{x}_0^\top \boldsymbol{\beta}$  can therefore be represented by the following two “observations”:

$$(21.0.30) \quad \begin{bmatrix} \hat{y}_0 \\ \mathbf{y}_0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mathbf{x}_0^\top \boldsymbol{\beta} + \begin{bmatrix} \delta_1 \\ \delta_2 \end{bmatrix} \quad \begin{bmatrix} \delta_1 \\ \delta_2 \end{bmatrix} \sim \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma^2 \begin{bmatrix} \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

This is a regression model with two observations and one unknown parameter,  $\mathbf{x}_0^\top \boldsymbol{\beta}$ , which has a nonspherical error covariance matrix. The formula for the BLUE of  $\mathbf{x}_0^\top \boldsymbol{\beta}$  in model (21.0.30) is

$$(21.0.31) \quad \hat{y}_0 = \left( [1 \quad 1] \begin{bmatrix} \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 & 0 \\ 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^{-1} [1 \quad 1] \begin{bmatrix} \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 & 0 \\ 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} \hat{y}_0 \\ \mathbf{y}_0 \end{bmatrix}$$

$$(21.0.32) \quad = \frac{1}{1 + \frac{1}{\mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0}} \left( \frac{\hat{y}_0}{\mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0} + \mathbf{y}_0 \right)$$

$$(21.0.33) \quad = \frac{1}{1 + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0} (\hat{y}_0 + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 \mathbf{y}_0).$$

Now subtract (21.0.33) from  $\mathbf{y}_0$  to get (21.0.29). □

Using (21.0.29), one can write (21.0.25) as

$$(21.0.34) \quad \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 \hat{\varepsilon}_0$$

Later, in (25.4.1), one will see that it can also be written in the form

$$(21.0.35) \quad \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} + (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{x}_0 (\mathbf{y}_0 - \mathbf{x}_0^\top \hat{\boldsymbol{\beta}})$$

where  $\mathbf{Z} = \begin{bmatrix} \mathbf{X} \\ \mathbf{x}_0^\top \end{bmatrix}$ .

**PROBLEM 279.** Show the following fact which is point (5) in the above updating principle in this special case: If one takes the squares of the standardized predictive residuals, one gets the difference of the SSE for the regression with and without the additional observation  $\mathbf{y}_0$

$$(21.0.36) \quad \text{SSE}^* - \text{SSE} = \frac{(\mathbf{y}_0 - \mathbf{x}_0^\top \hat{\boldsymbol{\beta}})^2}{1 + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0}$$

**ANSWER.** The sum of squared errors in the old regression is  $\text{SSE} = (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})$ ; the sum of squared errors in the updated regression is  $\text{SSE}^* = (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) + \hat{\varepsilon}_0^2$ . From (21.0.34) follows

$$(21.0.37) \quad \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 \hat{\varepsilon}_0.$$

If one squares this, the cross product terms fall away:  $(\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) = (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) + \hat{\varepsilon}_0 \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0 \hat{\varepsilon}_0$ . Adding  $\hat{\varepsilon}_0^2$  to both sides gives  $\text{SSE}^* = \text{SSE} + \hat{\varepsilon}_0^2 (1 + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0)$ . Now use (21.0.29) to get (21.0.36). □

## CHAPTER 22

## Constrained Least Squares

One of the assumptions for the linear model was that nothing is known about the true value of  $\beta$ . Any  $k$ -vector  $\gamma$  is a possible candidate for the value of  $\beta$ . We used this assumption e.g. when we concluded that an unbiased estimator  $\hat{\mathbf{B}}\mathbf{y}$  of  $\beta$  must satisfy  $\hat{\mathbf{B}}\mathbf{X} = \mathbf{I}$ . Now we will modify this assumption and assume we know that the true value  $\beta$  satisfies the linear constraint  $\mathbf{R}\beta = \mathbf{u}$ . To fix notation, assume  $\mathbf{y}$  be a  $n \times 1$  vector,  $\mathbf{u}$  a  $i \times 1$  vector,  $\mathbf{X}$  a  $n \times k$  matrix, and  $\mathbf{R}$  a  $i \times k$  matrix. In addition to our usual assumption that all columns of  $\mathbf{X}$  are linearly independent (i.e.,  $\mathbf{X}$  has full column rank) we will also make the assumption that all rows of  $\mathbf{R}$  are linearly independent (which is called:  $\mathbf{R}$  has full row rank). In other words, the matrix of constraints  $\mathbf{R}$  does not include “redundant” constraints which are linear combinations of the other constraints.

## 22.1. Building the Constraint into the Model

PROBLEM 280. Given a regression with a constant term and two explanatory variables which we will call  $x$  and  $z$ , i.e.,

$$(22.1.1) \quad \mathbf{y}_t = \alpha + \beta x_t + \gamma z_t + \varepsilon_t$$

- a. 1 point How will you estimate  $\beta$  and  $\gamma$  if it is known that  $\beta = \gamma$ ?

ANSWER. Write

$$(22.1.2) \quad \mathbf{y}_t = \alpha + \beta(x_t + z_t) + \varepsilon_t$$

- b. 1 point How will you estimate  $\beta$  and  $\gamma$  if it is known that  $\beta + \gamma = 1$ ?

ANSWER. Setting  $\gamma = 1 - \beta$  gives the regression

$$(22.1.3) \quad \mathbf{y}_t - z_t = \alpha + \beta(x_t - z_t) + \varepsilon_t$$

- c. 3 points Go back to a. If you add the original  $z$  as an additional regressor into the modified regression incorporating the constraint  $\beta = \gamma$ , then the coefficient of  $z$  is no longer an estimate of the original  $\gamma$ , but of a new parameter  $\delta$  which is a

linear combination of  $\alpha$ ,  $\beta$ , and  $\gamma$ . Compute this linear combination, i.e., express in terms of  $\alpha$ ,  $\beta$ , and  $\gamma$ . Remark (no proof required): this regression is equivalent (22.1.1), and it allows you to test the constraint.

ANSWER. If you add  $z$  as additional regressor into (22.1.2), you get  $\mathbf{y}_t = \alpha + \beta(x_t + z_t) + \delta z_t + \varepsilon_t$ . Now substitute the right hand side from (22.1.1) for  $\mathbf{y}$  to get  $\alpha + \beta x_t + \gamma z_t + \varepsilon_t = \alpha + \beta(x_t + z_t) + \delta z_t + \varepsilon_t$ . Cancelling out gives  $\gamma z_t = \beta z_t + \delta z_t$ , in other words,  $\gamma = \beta + \delta$ . In this regression therefore, the coefficient of  $z$  is split into the sum of two terms, the first term is the value it should be if the constraint were satisfied, and the other term is the difference from that.

- d. 2 points Now do the same thing with the modified regression from part b which incorporates the constraint  $\beta + \gamma = 1$ : include the original  $z$  as an additional regressor and determine the meaning of the coefficient of  $z$ .

What Problem 280 suggests is true in general: every constrained Least Squares problem can be reduced to an equivalent unconstrained Least Squares problem with fewer explanatory variables. Indeed, one can consider every least squares problem to be “constrained” because the assumption  $\mathcal{E}[\mathbf{y}] = \mathbf{X}\beta$  for some  $\beta$  is equivalent to a linear constraint on  $\mathcal{E}[\mathbf{y}]$ . The decision not to include certain explanatory variables in the regression can be considered the decision to set certain elements of  $\beta$  zero, which is the imposition of a constraint. If one writes a certain regression model as a constrained version of some other regression model, this simply means that one is interested in the relationship between two nested regressions.

Problem 219 is another example here.

## 22.2. Conversion of an Arbitrary Constraint into a Zero Constraint

This section, which is nothing but the matrix version of Problem 280, follows [DM93, pp. 16–19]. By reordering the elements of  $\beta$  one can write the constraint  $\mathbf{R}\beta = \mathbf{u}$  in the form

$$(22.2.1) \quad [\mathbf{R}_1 \quad \mathbf{R}_2] \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \equiv \mathbf{R}_1\beta_1 + \mathbf{R}_2\beta_2 = \mathbf{u}$$

where  $\mathbf{R}_1$  is a nonsingular  $i \times i$  matrix. Why can that be done? The rank of  $\mathbf{R}$  is  $i$ , i.e., all the rows are linearly independent. Since row rank is equal to column rank, there are also  $i$  linearly independent columns. Use those for  $\mathbf{R}_1$ . Using this subpartition, the original regression can be written

$$(22.2.2) \quad \mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \varepsilon$$

Now one can solve (22.2.1) for  $\beta_1$  to get

$$(22.2.3) \quad \beta_1 = \mathbf{R}_1^{-1}\mathbf{u} - \mathbf{R}_1^{-1}\mathbf{R}_2\beta_2$$

Plug (22.2.3) into (22.2.2) and rearrange to get a regression which is equivalent to the constrained regression:

$$(22.2.4) \quad \mathbf{y} - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{u} = (\mathbf{X}_2 - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{R}_2) \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

or

$$(22.2.5) \quad \mathbf{y}^* = \mathbf{Z}_2 \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

One more thing is noteworthy here: if we add  $\mathbf{X}_1$  as additional regressors into (22.2.5), we get a regression that is equivalent to (22.2.2). To see this, define the difference between the left hand side and right hand side of (22.2.3) as  $\boldsymbol{\gamma}_1 = \boldsymbol{\beta}_1 - \mathbf{R}_1^{-1} \mathbf{u} + \mathbf{R}_1^{-1} \mathbf{R}_2 \boldsymbol{\beta}_2$ ; then the constraint (22.2.1) is equivalent to the “zero constraint”  $\boldsymbol{\gamma}_1 = \mathbf{o}$ , and the regression

$$(22.2.6) \quad \mathbf{y} - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{u} = (\mathbf{X}_2 - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{R}_2) \boldsymbol{\beta}_2 + \mathbf{X}_1 (\boldsymbol{\beta}_1 - \mathbf{R}_1^{-1} \mathbf{u} + \mathbf{R}_1^{-1} \mathbf{R}_2 \boldsymbol{\beta}_2) + \boldsymbol{\varepsilon}$$

is equivalent to the original regression (22.2.2). (22.2.6) can also be written as

$$(22.2.7) \quad \mathbf{y}^* = \mathbf{Z}_2 \boldsymbol{\beta}_2 + \mathbf{X}_1 \boldsymbol{\gamma}_1 + \boldsymbol{\varepsilon}$$

The coefficient of  $\mathbf{X}_1$ , if it is added back into (22.2.5), is therefore  $\boldsymbol{\gamma}_1$ .

PROBLEM 281. [DM93] assert on p. 17, middle, that

$$(22.2.8) \quad \mathbf{R}[\mathbf{X}_1, \mathbf{Z}_2] = \mathbf{R}[\mathbf{X}_1, \mathbf{X}_2].$$

where  $\mathbf{Z}_2 = \mathbf{X}_2 - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{R}_2$ . Give a proof.

ANSWER. We have to show

$$(22.2.9) \quad \{z: z = \mathbf{X}_1 \boldsymbol{\gamma} + \mathbf{X}_2 \boldsymbol{\delta}\} = \{z: z = \mathbf{X}_1 \boldsymbol{\alpha} + \mathbf{Z}_2 \boldsymbol{\beta}\}$$

First  $\subset$ : given  $\boldsymbol{\gamma}$  and  $\boldsymbol{\delta}$  we need a  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  with

$$(22.2.10) \quad \mathbf{X}_1 \boldsymbol{\gamma} + \mathbf{X}_2 \boldsymbol{\delta} = \mathbf{X}_1 \boldsymbol{\alpha} + (\mathbf{X}_2 - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{R}_2) \boldsymbol{\beta}$$

This can be accomplished with  $\boldsymbol{\beta} = \boldsymbol{\delta}$  and  $\boldsymbol{\alpha} = \boldsymbol{\gamma} + \mathbf{R}_1^{-1} \mathbf{R}_2 \boldsymbol{\delta}$ . The other side is even more trivial: given  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$ , multiplying out the right side of (22.2.10) gives  $\mathbf{X}_1 \boldsymbol{\alpha} + \mathbf{X}_2 \boldsymbol{\beta} - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{R}_2 \boldsymbol{\beta}$ , i.e.,  $\boldsymbol{\delta} = \boldsymbol{\beta}$  and  $\boldsymbol{\gamma} = \boldsymbol{\alpha} - \mathbf{R}_1^{-1} \mathbf{R}_2 \boldsymbol{\beta}$ .  $\square$

### 22.3. Lagrange Approach to Constrained Least Squares

The constrained least squares estimator is that  $k \times 1$  vector  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$  which minimizes  $\text{SSE} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  subject to the linear constraint  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$ .

Again, we assume that  $\mathbf{X}$  has full column and  $\mathbf{R}$  full row rank.

The Lagrange approach to constrained least squares, which we follow here, is given in [Gre97, Section 7.3 on pp. 341/2], also [DM93, pp. 90/1]:

The Constrained Least Squares problem can be solved with the help of the “Lagrange function,” which is a function of the  $k \times 1$  vector  $\boldsymbol{\beta}$  and an additional  $i$  vector  $\boldsymbol{\lambda}$  of “Lagrange multipliers”:

$$(22.3.1) \quad \mathbf{L}(\boldsymbol{\beta}, \boldsymbol{\lambda}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\mathbf{R}\boldsymbol{\beta} - \mathbf{u})^\top \boldsymbol{\lambda}$$

$\boldsymbol{\lambda}$  can be considered a vector of “penalties” for violating the constraint. For every possible value of  $\boldsymbol{\lambda}$  one computes that  $\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}$  which minimizes  $\mathbf{L}$  for that  $\boldsymbol{\lambda}$  (This is an unconstrained minimization problem.) It will turn out that for one of the values  $\boldsymbol{\lambda} = \boldsymbol{\lambda}^*$ , the corresponding  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$  satisfies the constraint. This  $\hat{\boldsymbol{\beta}}$  is the solution to the constrained minimization problem we are looking for.

PROBLEM 282. 4 points Show the following: If  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$  is the unconstrained minimum argument of the Lagrange function

$$(22.3.2) \quad \mathbf{L}(\boldsymbol{\beta}, \boldsymbol{\lambda}^*) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\mathbf{R}\boldsymbol{\beta} - \mathbf{u})^\top \boldsymbol{\lambda}^*$$

for some fixed value  $\boldsymbol{\lambda}^*$ , and if at the same time  $\hat{\boldsymbol{\beta}}$  satisfies  $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{u}$ , then  $\hat{\boldsymbol{\beta}}$  minimizes  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  subject to the constraint  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$ .

ANSWER. Since  $\hat{\boldsymbol{\beta}}$  minimizes the Lagrange function, we know that

$$(22.3.3) \quad (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})^\top \boldsymbol{\lambda}^* \geq (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) + (\mathbf{R}\tilde{\boldsymbol{\beta}} - \mathbf{u})^\top \boldsymbol{\lambda}^*$$

for all  $\tilde{\boldsymbol{\beta}}$ . Since by assumption,  $\hat{\boldsymbol{\beta}}$  also satisfies the constraint, this simplifies to:

$$(22.3.4) \quad (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) + (\mathbf{R}\tilde{\boldsymbol{\beta}} - \mathbf{u})^\top \boldsymbol{\lambda}^* \geq (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

This is still true for all  $\tilde{\boldsymbol{\beta}}$ . If we only look at those  $\tilde{\boldsymbol{\beta}}$  which satisfy the constraint, we get

$$(22.3.5) \quad (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) \geq (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

This means,  $\hat{\boldsymbol{\beta}}$  is the constrained minimum argument.

Instead of imposing the constraint itself, one imposes a penalty function which has such a form that the agents will “voluntarily” heed the constraint. This is a familiar principle in neoclassical economics: instead of restricting pollution to a certain level, tax the polluters so much that they will voluntarily stay within the desired level.

The proof which follows now not only derives the formula for  $\hat{\boldsymbol{\beta}}$  but also shows that there is always a  $\boldsymbol{\lambda}^*$  for which  $\hat{\boldsymbol{\beta}}$  satisfies  $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{u}$ .

PROBLEM 283. 2 points Use the simple matrix differentiation rules  $\partial(\mathbf{w}^\top \boldsymbol{\beta})/\partial \boldsymbol{\beta} = \mathbf{w}^\top$  and  $\partial(\boldsymbol{\beta}^\top \mathbf{M}\boldsymbol{\beta})/\partial \boldsymbol{\beta}^\top = 2\boldsymbol{\beta}^\top \mathbf{M}$  to compute  $\partial \mathbf{L}/\partial \boldsymbol{\beta}^\top$  where

$$(22.3.6) \quad \mathbf{L}(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\mathbf{R}\boldsymbol{\beta} - \mathbf{u})^\top \boldsymbol{\lambda}$$

ANSWER. Write the objective function as  $\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\lambda}^\top \mathbf{R}\boldsymbol{\beta} - \boldsymbol{\lambda}^\top \mathbf{u}$  to get (22.3.7).  $\square$

Our goal is to find a  $\hat{\boldsymbol{\beta}}$  and a  $\boldsymbol{\lambda}^*$  so that (a)  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$  minimizes  $L(\boldsymbol{\beta}, \boldsymbol{\lambda}^*)$  and (b)  $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{u}$ . In other words,  $\hat{\boldsymbol{\beta}}$  and  $\boldsymbol{\lambda}^*$  together satisfy the following two conditions: (a) they must satisfy the first order condition for the unconstrained minimization of  $L$  with respect to  $\boldsymbol{\beta}$ , i.e.,  $\hat{\boldsymbol{\beta}}$  must annul

$$(22.3.7) \quad \partial L / \partial \boldsymbol{\beta}^\top = -2\mathbf{y}^\top \mathbf{X} + 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} + \boldsymbol{\lambda}^{*\top} \mathbf{R},$$

and (b)  $\hat{\boldsymbol{\beta}}$  must satisfy the constraint (22.3.9).

(22.3.7) and (22.3.9) are two linear matrix equations which can indeed be solved for  $\hat{\boldsymbol{\beta}}$  and  $\boldsymbol{\lambda}^*$ . I wrote (22.3.7) as a row vector, because the Jacobian of a scalar function is a row vector, but it is usually written as a column vector. Since this conventional notation is arithmetically a little simpler here, we will replace (22.3.7) with its transpose (22.3.8). Our starting point is therefore

$$(22.3.8) \quad 2\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = 2\mathbf{X}^\top \mathbf{y} - \mathbf{R}^\top \boldsymbol{\lambda}^*$$

$$(22.3.9) \quad \mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u} = \mathbf{o}$$

Some textbook treatments have an extra factor 2 in front of  $\boldsymbol{\lambda}^*$ , which makes the math slightly smoother, but which has the disadvantage that the Lagrange multiplier can no longer be interpreted as the “shadow price” for violating the constraint.

Solve (22.3.8) for  $\hat{\boldsymbol{\beta}}$  to get that  $\hat{\boldsymbol{\beta}}$  which minimizes  $L$  for any given  $\boldsymbol{\lambda}^*$ :

$$(22.3.10) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \frac{1}{2} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}^* = \hat{\boldsymbol{\beta}} - \frac{1}{2} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}^*$$

Here  $\hat{\boldsymbol{\beta}}$  on the right hand side is the *unconstrained* OLS estimate. Plug this formula for  $\hat{\boldsymbol{\beta}}$  into (22.3.9) in order to determine that value of  $\boldsymbol{\lambda}^*$  for which the corresponding  $\hat{\boldsymbol{\beta}}$  satisfies the constraint:

$$(22.3.11) \quad \mathbf{R}\hat{\boldsymbol{\beta}} - \frac{1}{2} \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}^* - \mathbf{u} = \mathbf{o}.$$

Since  $\mathbf{R}$  has full row rank and  $\mathbf{X}$  full column rank,  $\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top$  has an inverse (Problem 284). Therefore one can solve for  $\boldsymbol{\lambda}^*$ :

$$(22.3.12) \quad \boldsymbol{\lambda}^* = 2(\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})$$

If one substitutes this  $\boldsymbol{\lambda}^*$  back into (22.3.10), one gets the formula for the constrained least squares estimator:

$$(22.3.13) \quad \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}).$$

PROBLEM 284. If  $\mathbf{R}$  has full row rank and  $\mathbf{X}$  full column rank, show that  $\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top$  has an inverse.

ANSWER. Since it is nonnegative definite we have to show that it is positive definite.  $\mathbf{b}^\top \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \mathbf{b}$  implies  $\mathbf{b}^\top \mathbf{R} = \mathbf{o}^\top$  because  $(\mathbf{X}^\top \mathbf{X})^{-1}$  is positive definite, and this implies  $\mathbf{b} = \mathbf{o}$  because  $\mathbf{R}$  has full row rank.

PROBLEM 285. Assume  $\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2 \boldsymbol{\Psi})$  with a nonsingular  $\boldsymbol{\Psi}$  and show: If  $\hat{\boldsymbol{\beta}}$  minimizes SSE =  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  subject to the linear constraint  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$  the formula for the minimum argument  $\hat{\boldsymbol{\beta}}$  is the following modification of (22.3.13)

$$(22.3.14) \quad \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})$$

where  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{y}$ . This formula is given in [JHG+88, (11.2.1) on p. 457]. Remark, which you are not asked to prove: this is the best linear unbiased estimator if  $\boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2 \boldsymbol{\Psi})$  among all linear estimators which are unbiased whenever the true  $\boldsymbol{\beta}$  satisfies the constraint  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$ .

ANSWER. Lagrange function is

$$\begin{aligned} L(\boldsymbol{\beta}, \boldsymbol{\lambda}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\mathbf{R}\boldsymbol{\beta} - \mathbf{u})^\top \boldsymbol{\lambda} \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \boldsymbol{\Psi}^{-1} \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\lambda}^\top \mathbf{R}\boldsymbol{\beta} - \boldsymbol{\lambda}^\top \mathbf{u} \end{aligned}$$

Jacobian is

$$\partial L / \partial \boldsymbol{\beta}^\top = -2\mathbf{y}^\top \boldsymbol{\Psi}^{-1} \mathbf{X} + 2\boldsymbol{\beta}^\top \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X} + \boldsymbol{\lambda}^\top \mathbf{R},$$

Transposing and setting it zero gives

$$(22.3.15) \quad 2\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X} \hat{\boldsymbol{\beta}} = 2\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{y} - \mathbf{R}^\top \boldsymbol{\lambda}^*$$

Solve (22.3.15) for  $\hat{\boldsymbol{\beta}}$ :

$$(22.3.16) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{y} - \frac{1}{2} (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}^* = \hat{\boldsymbol{\beta}} - \frac{1}{2} (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}^*$$

Here  $\hat{\boldsymbol{\beta}}$  is the unconstrained GLS estimate. Plug  $\hat{\boldsymbol{\beta}}$  into the constraint (22.3.9):

$$(22.3.17) \quad \mathbf{R}\hat{\boldsymbol{\beta}} - \frac{1}{2} \mathbf{R}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}^* - \mathbf{u} = \mathbf{o}.$$

Since  $\mathbf{R}$  has full row rank and  $\mathbf{X}$  full column rank and  $\boldsymbol{\Psi}$  is nonsingular,  $\mathbf{R}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top$  has an inverse. Therefore

$$(22.3.18) \quad \boldsymbol{\lambda}^* = 2(\mathbf{R}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})$$

Now substitute this  $\boldsymbol{\lambda}^*$  back into (22.3.16):

$$(22.3.19) \quad \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \boldsymbol{\Psi}^{-1} \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}).$$

### 22.4. Constrained Least Squares as the Nesting of Two Simpler Models

The imposition of a constraint can also be considered the addition of new information: a certain linear transformation of  $\beta$ , namely,  $R\beta$ , is observed without error.

PROBLEM 286. Assume the random  $\beta \sim (\hat{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$  is unobserved, but one observes  $R\beta = \mathbf{u}$ .

- a. 2 points Compute the best linear predictor of  $\beta$  on the basis of the observation  $\mathbf{u}$ . Hint: First write down the joint means and covariance matrix of  $\mathbf{u}$  and  $\beta$ .

ANSWER.

$$(22.4.1) \quad \begin{bmatrix} \mathbf{u} \\ \beta \end{bmatrix} \sim \left( \begin{bmatrix} R\hat{\beta} \\ \hat{\beta} \end{bmatrix}, \sigma^2 \begin{bmatrix} R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top & R(\mathbf{X}^\top \mathbf{X})^{-1} \\ (\mathbf{X}^\top \mathbf{X})^{-1}R^\top & (\mathbf{X}^\top \mathbf{X})^{-1} \end{bmatrix} \right).$$

Therefore application of formula (??) gives

$$(22.4.2) \quad \beta^* = \hat{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1}R^\top (R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top)^{-1}(\mathbf{u} - R\hat{\beta}).$$

□

- b. 1 point Look at the formula for the predictor you just derived. Have you seen this formula before? Describe the situation in which this formula is valid as a BLUE-formula, and compare the situation with the situation here.

ANSWER. Of course, constrained least squares. But in constrained least squares,  $\beta$  is nonrandom and  $\hat{\beta}$  is random, while here it is the other way round. □

In the unconstrained OLS model, i.e., before the “observation” of  $\mathbf{u} = R\beta$ , the best bounded  $MSE$  estimators of  $\mathbf{u}$  and  $\beta$  are  $R\hat{\beta}$  and  $\hat{\beta}$ , with the sampling errors having the following means and variances:

$$(22.4.3) \quad \begin{bmatrix} \mathbf{u} - R\hat{\beta} \\ \beta - \hat{\beta} \end{bmatrix} \sim \left( \begin{bmatrix} \mathbf{o} \\ \mathbf{o} \end{bmatrix}, \sigma^2 \begin{bmatrix} R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top & R(\mathbf{X}^\top \mathbf{X})^{-1} \\ (\mathbf{X}^\top \mathbf{X})^{-1}R^\top & (\mathbf{X}^\top \mathbf{X})^{-1} \end{bmatrix} \right)$$

After the observation of  $\mathbf{u}$  we can therefore apply (20.1.18) to get exactly equation (22.3.13) for  $\hat{\beta}$ . This is probably the easiest way to derive this equation, but it derives constrained least squares by the minimization of the  $MSE$ -matrix, not by the least squares problem.

### 22.5. Solution by Quadratic Decomposition

An alternative purely algebraic solution method for this constrained minimization problem rewrites the OLS objective function in such a way that one sees immediately what the constrained minimum value is.

Start with the decomposition (14.2.12) which can be used to show optimality of the OLS estimate:

$$(\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) = (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) + (\beta - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\beta - \hat{\beta}).$$

Split the second term again, using  $\hat{\beta} - \hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1}R^\top (R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top)^{-1}(R\hat{\beta} - \mathbf{u})$ :

$$\begin{aligned} (\beta - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\beta - \hat{\beta}) &= (\beta - \hat{\beta} - (\hat{\beta} - \hat{\beta}))^\top \mathbf{X}^\top \mathbf{X} (\beta - \hat{\beta} - (\hat{\beta} - \hat{\beta})) \\ &= (\beta - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\beta - \hat{\beta}) \\ &\quad - 2(\beta - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}R^\top (R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top)^{-1}(R\hat{\beta} - \mathbf{u}) \\ &\quad + (\hat{\beta} - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \hat{\beta}). \end{aligned}$$

The cross product terms can be simplified to  $-2(R\beta - \mathbf{u})^\top (R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top)^{-1}(R\hat{\beta} - \mathbf{u})$ , and the last term is  $(R\hat{\beta} - \mathbf{u})^\top (R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top)^{-1}(R\hat{\beta} - \mathbf{u})$ . Therefore the objective function for an arbitrary  $\beta$  can be written as

$$\begin{aligned} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) &= (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) \\ &\quad + (\beta - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\beta - \hat{\beta}) \\ &\quad - 2(R\beta - \mathbf{u})^\top (R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top)^{-1}(R\hat{\beta} - \mathbf{u}) \\ &\quad + (R\hat{\beta} - \mathbf{u})^\top (R(\mathbf{X}^\top \mathbf{X})^{-1}R^\top)^{-1}(R\hat{\beta} - \mathbf{u}) \end{aligned}$$

The first and last terms do not depend on  $\beta$  at all; the third term is zero whenever  $\beta$  satisfies  $R\beta = \mathbf{u}$ ; and the second term is minimized if and only if  $\beta = \hat{\beta}$ , in which case it also takes the value zero.

### 22.6. Sampling Properties of Constrained Least Squares

Again, this variant of the least squares principle leads to estimators with desirable sampling properties. Note that  $\hat{\beta}$  is an affine function of  $\mathbf{y}$ . We will compute  $\mathcal{E}[\hat{\beta} - \beta]$  and  $MSE[\hat{\beta}; \beta]$  not only in the case that the true  $\beta$  satisfies  $R\beta = \mathbf{u}$ , but also in the case that it does not. For this, let us first get a suitable representation of  $\hat{\beta}$ .

sampling error:

$$\begin{aligned}
 \hat{\beta} - \beta &= (\hat{\beta} - \beta) + (\hat{\beta} - \hat{\beta}) = \\
 (22.6.1) \quad &= (\hat{\beta} - \beta) - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \mathbf{R}(\hat{\beta} - \beta) \\
 &\quad - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\beta - \mathbf{u}).
 \end{aligned}$$

The last term is zero if  $\beta$  satisfies the constraint. Now use (18.0.7) twice to get

$$(22.6.2) \quad \hat{\beta} - \beta = \mathbf{W} \mathbf{X}^\top \varepsilon - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\beta - \mathbf{u})$$

where

$$(22.6.3) \quad \mathbf{W} = (\mathbf{X}^\top \mathbf{X})^{-1} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}.$$

If  $\beta$  satisfies the constraint, (22.6.2) simplifies to  $\hat{\beta} - \beta = \mathbf{W} \mathbf{X}^\top \varepsilon$ . In this case, therefore,  $\hat{\beta}$  is unbiased and  $\mathcal{MSE}[\hat{\beta}; \beta] = \sigma^2 \mathbf{W}$  (Problem 287). Since  $(\mathbf{X}^\top \mathbf{X})^{-1} - \mathbf{W}$  is nonnegative definite,  $\mathcal{MSE}[\hat{\beta}; \beta]$  is smaller than  $\mathcal{MSE}[\hat{\beta}; \beta]$  by a nonnegative definite matrix. This should be expected, since  $\hat{\beta}$  uses more information than  $\hat{\beta}$ .

PROBLEM 287.

- a. Show that  $\mathbf{W} \mathbf{X}^\top \mathbf{X} \mathbf{W} = \mathbf{W}$  (i.e.,  $\mathbf{X}^\top \mathbf{X}$  is a  $g$ -inverse of  $\mathbf{W}$ ).

ANSWER. This is a tedious matrix multiplication. □

- b. Use this to show that  $\mathcal{MSE}[\hat{\beta}; \beta] = \sigma^2 \mathbf{W}$ .

(Without proof:) The Gauss-Markov theorem can be extended here as follows: the constrained least squares estimator is the best linear unbiased estimator among all linear (or, more precisely, affine) estimators which are unbiased whenever the true  $\beta$  satisfies the constraint  $\mathbf{R}\beta = \mathbf{u}$ . Note that there are more estimators which are unbiased whenever the true  $\beta$  satisfies the constraint than there are estimators which are unbiased for all  $\beta$ .

If  $\mathbf{R}\beta \neq \mathbf{u}$ , then  $\hat{\beta}$  is biased. Its bias is

$$(22.6.4) \quad \mathcal{E}[\hat{\beta} - \beta] = -(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\beta - \mathbf{u}).$$

Due to the decomposition (17.1.2) of the  $\mathcal{MSE}$  matrix into dispersion matrix plus squared bias, it follows

$$\begin{aligned}
 (22.6.5) \quad \mathcal{MSE}[\hat{\beta}; \beta] &= \sigma^2 \mathbf{W} + \\
 &\quad + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\beta - \mathbf{u}) \\
 &\quad \cdot (\mathbf{R}\beta - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}
 \end{aligned}$$

Even if the true parameter does not satisfy the constraint, it is still possible that the constrained least squares estimator has a better  $\mathcal{MSE}$  matrix than the unconstrained one. This is the case if and only if the true parameter values  $\beta$  and  $\sigma^2$  satisfy

$$(22.6.6) \quad (\mathbf{R}\beta - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\beta - \mathbf{u}) \leq \sigma^2.$$

This equation, which is the same as [Gre97, (8-27) on p. 406], is an interesting result, because the obvious estimate of the lefthand side in (22.6.6) is  $i$  times the value of the F-test statistic for the hypothesis  $\mathbf{R}\beta = \mathbf{u}$ . To test for this, one has to use the *noncentral* F-test with parameters  $i$ ,  $n - k$ , and  $1/2$ .

PROBLEM 288. *2 points* This Problem motivates Equation (22.6.6). If  $\hat{\beta}$  is a better estimator of  $\beta$  than  $\hat{\beta}$ , then  $\mathbf{R}\hat{\beta} = \mathbf{u}$  is also a better estimator of  $\mathbf{R}\beta$  than  $\mathbf{R}\hat{\beta}$ . Show that this latter condition is not only necessary but already sufficient, i.e., if  $\mathcal{MSE}[\mathbf{R}\hat{\beta}; \mathbf{R}\beta] - \mathcal{MSE}[\mathbf{u}; \mathbf{R}\beta]$  is nonnegative definite then  $\beta$  and  $\sigma^2$  satisfy (22.6.6). You are allowed to use, without proof, theorem A.5.9 in the mathematical Appendix.

ANSWER. We have to show

$$(22.6.7) \quad \sigma^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top - (\mathbf{R}\beta - \mathbf{u})(\mathbf{R}\beta - \mathbf{u})^\top$$

is nonnegative definite. Since  $\mathbf{\Omega} = \sigma^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top$  has an inverse, theorem A.5.9 immediately leads to (22.6.6).

## 22.7. Estimation of the Variance in Constrained OLS

Next we will compute the expected value of the minimum value of the constrained OLS objective function, i.e.,  $\mathcal{E}[\hat{\varepsilon}^\top \hat{\varepsilon}]$  where  $\hat{\varepsilon} = \mathbf{y} - \mathbf{X}\hat{\beta}$ , again without necessarily making the assumption that  $\mathbf{R}\beta = \mathbf{u}$ :

$$(22.7.1) \quad \hat{\varepsilon} = \mathbf{y} - \mathbf{X}\hat{\beta} = \hat{\varepsilon} + \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\beta} - \mathbf{u}).$$

Since  $\mathbf{X}^\top \hat{\varepsilon} = \mathbf{o}$ , it follows

$$(22.7.2) \quad \hat{\varepsilon}^\top \hat{\varepsilon} = \hat{\varepsilon}^\top \hat{\varepsilon} + (\mathbf{R}\hat{\beta} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\beta} - \mathbf{u}).$$



Now note that  $\mathcal{E}[\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}] = \mathbf{R}\boldsymbol{\beta} - \mathbf{u}$  and  $\mathcal{V}[\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}] = \sigma^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top$ . Therefore use (??) in theorem ?? and  $\text{tr}\left((\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)(\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1}\right) = i$  to get

$$(22.7.3) \quad \begin{aligned} \mathbb{E}[(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})] &= \\ &= \sigma^2 i + (\mathbf{R}\boldsymbol{\beta} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\boldsymbol{\beta} - \mathbf{u}) \end{aligned}$$

Since  $\mathbb{E}[\hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}}] = \sigma^2(n - k)$ , it follows

$$(22.7.4) \quad \mathbb{E}[\hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}}] = \sigma^2(n + i - k) + (\mathbf{R}\boldsymbol{\beta} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\boldsymbol{\beta} - \mathbf{u}).$$

In other words,  $\hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}} / (n + i - k)$  is an unbiased estimator of  $\sigma^2$  if the constraint holds, and it is biased upwards if the constraint does not hold. The adjustment of the degrees of freedom is what one should expect: a regression with  $k$  explanatory variables and  $i$  constraints can always be rewritten as a regression with  $k - i$  different explanatory variables (see Section 22.2), and the distribution of the SSE does not depend on the values taken by the explanatory variables at all, only on how many there are. The unbiased estimate of  $\sigma^2$  is therefore

$$(22.7.5) \quad \hat{\sigma}^2 = \hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}} / (n + i - k)$$

Here is some geometric intuition:  $\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\epsilon}}$  is an orthogonal decomposition, since  $\hat{\boldsymbol{\epsilon}}$  is orthogonal to all columns of  $\mathbf{X}$ . From orthogonality follows  $\mathbf{y}^\top \mathbf{y} = \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}}$ . If one splits up  $\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\epsilon}}$ , one should expect this to be orthogonal as well. But this is only the case if  $\mathbf{u} = \mathbf{o}$ . If  $\mathbf{u} \neq \mathbf{o}$ , one first has to shift the origin of the coordinate system to a point which can be written in the form  $\mathbf{X}\boldsymbol{\beta}_0$  where  $\boldsymbol{\beta}_0$  satisfies the constraint:

PROBLEM 289. 3 points Assume  $\hat{\boldsymbol{\beta}}$  is the constrained least squares estimate, and  $\boldsymbol{\beta}_0$  is any vector satisfying  $\mathbf{R}\boldsymbol{\beta}_0 = \mathbf{u}$ . Show that in the decomposition

$$(22.7.6) \quad \mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0 = \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) + \hat{\boldsymbol{\epsilon}}$$

the two vectors on the righthand side are orthogonal.

ANSWER. We have to show  $(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = 0$ . Since  $\hat{\boldsymbol{\epsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) = \hat{\boldsymbol{\epsilon}} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ , and we already know that  $\mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = \mathbf{o}$ , it is necessary and sufficient to show that  $(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{X}^\top \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) = 0$ . By (22.3.13),

$$\begin{aligned} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{X}^\top \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) &= (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}) \\ &= (\mathbf{u} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}) = 0. \end{aligned}$$

□

If  $\mathbf{u} = \mathbf{o}$ , then one has two orthogonal decompositions:  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\boldsymbol{\epsilon}}$ , and  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\boldsymbol{\epsilon}}$ . And if one connects the footpoints of these two orthogonal decompositions, one obtains an orthogonal decomposition into three parts:

PROBLEM 290. Assume  $\hat{\boldsymbol{\beta}}$  is the constrained least squares estimator subject to the constraint  $\mathbf{R}\boldsymbol{\beta} = \mathbf{o}$ , and  $\boldsymbol{\beta}$  is the unconstrained least squares estimator.

- a. 1 point With the usual notation  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  and  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ , show that

$$(22.7.7) \quad \mathbf{y} = \hat{\mathbf{y}} + (\hat{\mathbf{y}} - \hat{\mathbf{y}}) + \hat{\boldsymbol{\epsilon}}$$

Point out these vectors in the *reggeom* simulation.

ANSWER. In the *reggeom*-simulation,  $\mathbf{y}$  is the purple line;  $\mathbf{X}\hat{\boldsymbol{\beta}}$  is the red line starting at origin, one could also call it  $\hat{\mathbf{y}}$ ;  $\mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) = \hat{\mathbf{y}} - \hat{\mathbf{y}}$  is the light blue line, and  $\hat{\boldsymbol{\epsilon}}$  is the green line which does not start at the origin. In other words: if one projects  $\mathbf{y}$  on a plane, and also on a line in that plane, and then connects the footpoints of these two projections, one obtains a zig-zag with two right angles.

- b. 4 points Show that in (22.7.7) the three vectors  $\hat{\mathbf{y}}$ ,  $\hat{\mathbf{y}} - \hat{\mathbf{y}}$ , and  $\hat{\boldsymbol{\epsilon}}$  are orthogonal. You are allowed to use, without proof, formula (22.3.13):

ANSWER. One has to verify that the scalar products of the three vectors on the right hand side of (22.7.7) are zero.  $\hat{\mathbf{y}}^\top \hat{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = 0$  and  $(\hat{\mathbf{y}} - \hat{\mathbf{y}})^\top \hat{\boldsymbol{\epsilon}} = (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = 0$  follow from  $\mathbf{X}^\top \hat{\boldsymbol{\epsilon}} = \mathbf{o}$ ; geometrically one can simply say that  $\hat{\mathbf{y}}$  and  $\hat{\mathbf{y}}$  are in the space spanned by the columns of  $\mathbf{X}$ , and  $\hat{\boldsymbol{\epsilon}}$  is orthogonal to that space. Finally, using (22.3.13) for  $\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}$ ,

$$\begin{aligned} \hat{\mathbf{y}}^\top (\hat{\mathbf{y}} - \hat{\mathbf{y}}) &= \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) = \\ &= \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \mathbf{R}\hat{\boldsymbol{\beta}} = \\ &= \hat{\boldsymbol{\beta}}^\top \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \mathbf{R}\hat{\boldsymbol{\beta}} = 0 \end{aligned}$$

because  $\hat{\boldsymbol{\beta}}$  satisfies the constraint  $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{o}$ , hence  $\hat{\boldsymbol{\beta}}^\top \mathbf{R}^\top = \mathbf{o}^\top$ .

PROBLEM 291.

- a. 3 points In the model  $\mathbf{y} = \boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\mathbf{y}$  is a  $n \times 1$  vector, and  $\boldsymbol{\epsilon} \sim (\mathbf{o}, \sigma^2)$  subject to the constraint  $\boldsymbol{\iota}^\top \boldsymbol{\beta} = 0$ , compute  $\hat{\boldsymbol{\beta}}$ ,  $\hat{\boldsymbol{\epsilon}}$ , and the unbiased estimate  $\hat{\sigma}^2$ . Give general formulas and the numerical results for the case  $\mathbf{y}^\top = [-1 \ 0 \ 1 \ 2]$ . You need to do is evaluate the appropriate formulas and correctly count the number of degrees of freedom.

ANSWER. The unconstrained least squares estimate of  $\boldsymbol{\beta}$  is  $\hat{\boldsymbol{\beta}} = \mathbf{y}$ , and since  $\mathbf{X} = \mathbf{I}$ ,  $\mathbf{R} = \boldsymbol{\iota}^\top$  and  $\mathbf{u} = 0$ , the constrained LSE has the form  $\hat{\boldsymbol{\beta}} = \mathbf{y} - \boldsymbol{\iota}(\boldsymbol{\iota}^\top \boldsymbol{\iota})^{-1}(\boldsymbol{\iota}^\top \mathbf{y}) = \mathbf{y} - \boldsymbol{\iota}\bar{y}$  by (22.3.13).  $\mathbf{y}^\top = [-1, 0, 1, 2]$  this gives  $\hat{\boldsymbol{\beta}}^\top = [-1.5, -0.5, 0.5, 1.5]$ . The residuals in the constrained model

therefore  $\hat{\epsilon} = \epsilon \bar{y}$ , i.e.,  $\hat{\epsilon} = [0.5, 0.5, 0.5, 0.5]$ . Since one has  $n$  observations,  $n$  parameters and 1 constraint, the number of degrees of freedom is 1. Therefore  $\hat{\sigma}^2 = \hat{\epsilon}^\top \hat{\epsilon} / 1 = n \bar{y}^2$  which is = 1 in our case.  $\square$

• b. 1 point Can you think of a practical situation in which this model might be appropriate?

ANSWER. This can occur if one measures data which theoretically add to zero, and the measurement errors are independent and have equal standard deviations.  $\square$

• c. 2 points Check your results against a SAS printout (or do it in any other statistical package) with the data vector  $\mathbf{y}^\top = [-1 \ 0 \ 1 \ 2]$ . Here are the sas commands:

```
data zeromean;
input y x1 x2 x3 x4;
cards;
-1 1 0 0 0
 0 0 1 0 0
 1 0 0 1 0
 2 0 0 0 1
;
proc reg;
model y= x1 x2 x3 x4 /
noint;
restrict x1+x2+x3+x4=0;
output out=zerout
residual=ehat;
run;
proc print data=zerout;
run;
```

PROBLEM 292. Least squares estimates of the coefficients of a linear regression model often have signs that are regarded by the researcher to be ‘wrong’. In an effort to obtain the ‘right’ signs, the researcher may be tempted to drop statistically insignificant variables from the equation. [Lea75] showed that such attempts necessarily fail: there can be no change in sign of any coefficient which is more significant than the coefficient of the omitted variable. The present exercise shows this, using a different proof than Leamer’s. You will need the formula for the constrained least squares estimator subject to one linear constraint  $\mathbf{r}^\top \boldsymbol{\beta} = u$ , which is

$$(22.7.8) \quad \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - \mathbf{V} \mathbf{r} (\mathbf{r}^\top \mathbf{V} \mathbf{r})^{-1} (\mathbf{r}^\top \hat{\boldsymbol{\beta}} - u).$$

where  $\mathbf{V} = (\mathbf{X}^\top \mathbf{X})^{-1}$ .

• a. In order to assess the sensitivity of the estimate of any linear combination of the elements of  $\boldsymbol{\beta}$ ,  $\phi = \mathbf{t}^\top \boldsymbol{\beta}$ , due to imposition of the constraint, it makes sense to divide the change  $\mathbf{t}^\top \hat{\boldsymbol{\beta}} - \mathbf{t}^\top \hat{\boldsymbol{\beta}}$  by the standard deviation of  $\mathbf{t}^\top \hat{\boldsymbol{\beta}}$ , i.e., to look at

$$(22.7.9) \quad \frac{\mathbf{t}^\top (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})}{\sigma \sqrt{\mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t}}}.$$

Such a standardization allows you to compare the sensitivity of different linear combinations. Show that that linear combination of the elements of  $\hat{\boldsymbol{\beta}}$  which is affected most if one imposes the constraint  $\mathbf{r}^\top \boldsymbol{\beta} = u$  is the constraint  $\mathbf{t} = \mathbf{r}$  itself. If  $\mathbf{t}^\top \hat{\boldsymbol{\beta}}$  value is small, then no other linear combination of the elements of  $\hat{\boldsymbol{\beta}}$  will be affected much by the imposition of the constraint either.

ANSWER. Using (22.7.8) and equation (25.4.1) one obtains

$$\begin{aligned} \max_{\mathbf{t}} \frac{(\mathbf{t}^\top (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}))^2}{\sigma^2 \mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t}} &= \frac{(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})}{\sigma^2} \\ &= \frac{(\mathbf{r}^\top \hat{\boldsymbol{\beta}} - u)^\top (\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r})^{-1} (\mathbf{r}^\top \hat{\boldsymbol{\beta}} - u)}{\sigma^2} = \frac{(\mathbf{r}^\top \hat{\boldsymbol{\beta}} - u)^2}{\sigma^2 \mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}} \end{aligned}$$

## 22.8. Inequality Restrictions

With linear inequality restrictions, it makes sense to have  $\mathbf{R}$  of deficient rank. These are like two different half planes in the same plane, and the restrictions define a quarter plane, or a triangle, etc.

One obvious approach would be: compute the unrestricted estimator, see which restrictions it violates, and apply these restrictions with equality. But this equal restricted estimator may then suddenly violate other restrictions.

One brute force approach would be: impose all combinations of restrictions and see if the so partially restricted parameter satisfies the other restrictions too; among those that do, choose the one with the lowest SSE.

[Gre97, 8.5.3 on pp. 411/12] has good discussion. The inequality restricted estimator is biased, unless the true parameter value satisfies all inequality restrictions with equality. It is always a mixture between the unbiased  $\hat{\boldsymbol{\beta}}$  and some restricted estimator which is biased if this condition does not hold.

Its variance is always smaller than that of  $\hat{\boldsymbol{\beta}}$  but, incredibly, its MSE will sometimes be larger than that of  $\hat{\boldsymbol{\beta}}$ . Don’t understand how this comes about.

### 22.9. Application: Biased Estimators and Pre-Test Estimators

The formulas about Constrained Least Squares which were just derived suggest that it is sometimes advantageous (in terms of MSE) to impose constraints even if they do not really hold. In other words, one should not put all explanatory variables into a regression which have an influence, but only the main ones. A logical extension of this idea is the common practice of first testing whether some variables have significant influence and dropping the variables if they do not. These so-called pre-test estimators are very common. [DM93, Chapter 3.7, pp. 94–98] says something about them. Pre-test estimation this seems a good procedure, but the graph regarding MSE shows it is not: the pre-test estimator never has lowest MSE, and it has highest MSE exactly in the area where it is most likely to be applied.

## CHAPTER 23

## Additional Regressors

A good detailed explanation of the topics covered in this chapter is [DM93, pp. 19–24]. [DM93] use the addition of variables as their main paradigm for going from a more restrictive to a less restrictive model.

In this chapter, the usual regression model is given in the form

$$(23.0.1) \quad \mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \boldsymbol{\varepsilon} = [\mathbf{X}_1 \quad \mathbf{X}_2] \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \boldsymbol{\varepsilon} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2\mathbf{I})$$

where  $\mathbf{X} = [\mathbf{X}_1 \quad \mathbf{X}_2]$  has full column rank, and the coefficient vector is  $\boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$ .

We take a sequential approach to this regression. First we regress  $\mathbf{y}$  on  $\mathbf{X}_1$  alone, which gives the regression coefficient  $\hat{\boldsymbol{\beta}}_1$ . This by itself is an inconsistent estimator of  $\beta_1$ , but we will use it as a stepping stone towards the full regression. We make use of the information gained by the regression on  $\mathbf{X}_1$  in our computation of the full regression. Such a sequential approach may be appropriate in the following situations:

- If regression on  $\mathbf{X}_1$  is much simpler than the combined regression, for instance if  $\mathbf{X}_1$  contains dummy or trend variables, and the dataset is large. Example: model (??).
- If we want to fit the regressors in  $\mathbf{X}_2$  by graphical methods and those in  $\mathbf{X}_1$  by analytical methods (added variable plots).
- If we need an estimate of  $\beta_2$  but are not interested in an estimate of  $\beta_1$ .
- If we want to test the joint significance of the regressors in  $\mathbf{X}_2$ , while  $\mathbf{X}_1$  consists of regressors not being tested.

If one regresses  $\mathbf{y}$  on  $\mathbf{X}_1$ , one gets  $\mathbf{y} = \mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \hat{\boldsymbol{\varepsilon}}$ . Of course,  $\hat{\boldsymbol{\beta}}_1$  is an inconsistent estimator of  $\beta_1$ , since some explanatory variables are left out. And  $\hat{\boldsymbol{\varepsilon}}$  is orthogonal to  $\mathbf{X}_1$  but not to  $\mathbf{X}_2$ .

The iterative “backfitting” method proceeds from here as follows: it regresses  $\hat{\boldsymbol{\varepsilon}}$  on  $\mathbf{X}_2$ , which gives another residual, which is again orthogonal on  $\mathbf{X}_2$  but no longer orthogonal on  $\mathbf{X}_1$ . Then this new residual is regressed on  $\mathbf{X}_1$  again, etc.

PROBLEM 293. *The purpose of this Problem is to get a graphical intuition of issues in sequential regression. Make sure the stand-alone program `xgobi` is installed on your computer (in Debian GNU-Linux do `apt-get install xgobi`), and the interface `xgobi` is installed (the R-command is simply `install.packages("xgobi")` or, on a Debian system the preferred argument is `install.packages("xgobi", lib = "/usr/lib/R/library")`). You have to give the commands `library(xgobi)` and then `reggeom()`. This produces a graph in the `XGobi` window which looks like [DM93, Figure 3b on p. 22]. If you switch from the `XYPlot` view to the `Rotation` view, you will see the same lines rotating 3-dimensionally, and you can interact with this graph. You will see that this graph shows the dependent variable  $\mathbf{y}$ , the regression of  $\mathbf{y}$  on  $\mathbf{x}_1$ , and the regression of  $\mathbf{y}$  on  $\mathbf{x}_1$  and  $\mathbf{x}_2$ .*

• a. 1 point In order to show that you have correctly identified which line is please answer the following two questions: Which color is  $\mathbf{y}$ : red, yellow, light blue, dark blue, green, purple, or white? If it is yellow, also answer the question: Is it the yellow line which is in part covered by a red line, or is it the other one? If it is red, green, or dark blue, also answer the question: Does it start at the origin or not?

• b. 1 point Now answer the same two questions about  $\mathbf{x}_1$ .

• c. 1 point Now answer the same two questions about  $\mathbf{x}_2$ .

• d. 1 point Now answer the same two questions about  $\hat{\boldsymbol{\varepsilon}}$ , the residual in regression of  $\mathbf{y}$  on  $\mathbf{x}_1$ .

• e. Now assume  $\mathbf{x}_1$  is the vector of ones. The  $R^2$  of this regression is a ratio of the squared lengths of two of the lines in the regression. Which lines?

• f. 2 points If one regresses  $\hat{\boldsymbol{\varepsilon}}$  on  $\mathbf{x}_2$ , one gets a decomposition  $\hat{\boldsymbol{\varepsilon}} = \mathbf{h} + \mathbf{k}$  where  $\mathbf{h}$  is a multiple of  $\mathbf{x}_2$  and  $\mathbf{k}$  orthogonal to  $\mathbf{x}_2$ . This is the next step in backfitting algorithm. Draw this decomposition into the diagram. The points are already invisibly present. Therefore you should use the line editor to connect the points. You may want to increase the magnification scale of the figure for this. (In my version of `XGobi`, I often lose lines if I try to add more lines. This seems to be a bug which will probably be fixed eventually.) Which label does the corner point in the decomposition have? Make a geometric argument that the new residual  $\mathbf{k}$  is no longer orthogonal to  $\mathbf{x}_2$ .

• g. 1 point The next step in the backfitting procedure is to regress  $\mathbf{k}$  on  $\mathbf{x}_1$ . The corner point for this decomposition is again invisibly in the animation. Identify the two endpoints of the residual in this regression. Hint: the R-command `example(reggeom)` produces a modified version of the animation in which the backfitting procedure is highlighted. The successive residuals which are used as regressors

are drawn in dark blue, and the quickly improving approximations to the fitted value are connected by a red zig-zag line.

• h. 1 point The diagram contains the points for two more backfitting steps. Identify the endpoints of both residuals.

• i. 2 points Of the five cornerpoints obtained by simple regressions,  $c, p, q, r,$  and  $s,$  three lie on one straight line, and the other two on a different straight line, with the intersection of these straight lines being the corner point in the multiple regression of  $\mathbf{y}$  on  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . Which three points are on the same line, and how can these two lines be characterized?

• j. 1 point Of the lines  $cp, pq, qr,$  and  $rs,$  two are parallel to  $\mathbf{x}_1,$  and two parallel to  $\mathbf{x}_2$ . Which two are parallel to  $\mathbf{x}_1$ ?

• k. 1 point Draw in the regression of  $\mathbf{y}$  on  $\mathbf{x}_2$ .

• l. 3 points Which two variables are plotted against each other in an added-variable plot for  $\mathbf{x}_2$ ?

Here are the coordinates of some of the points in this animation:

$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{y}$	$\hat{\mathbf{y}}$	$\hat{\hat{\mathbf{y}}}$
5	-1	3	3	3
0	4	3	3	0
0	0	4	0	0

In the dataset which R submits to XGobi, all coordinates are multiplied by 1156, which has the effect that all the points included in the animation have integer coordinates.

PROBLEM 294. 2 points How do you know that the decomposition  $\begin{bmatrix} 3 \\ 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 3 \\ 4 \end{bmatrix}$  is  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\hat{\mathbf{e}}}$  in the regression of  $\mathbf{y} = \begin{bmatrix} 3 \\ 3 \\ 4 \end{bmatrix}$  on  $\mathbf{x}_1 = \begin{bmatrix} 5 \\ 0 \\ 0 \end{bmatrix}$ ?

ANSWER. Besides the equation  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\hat{\mathbf{e}}}$  we have to check two things: (1)  $\hat{\mathbf{y}}$  is a linear combination of all the explanatory variables (here: is a multiple of  $\mathbf{x}_1$ ), and (2)  $\hat{\hat{\mathbf{e}}}$  is orthogonal to all explanatory variables. Compare Problem ??.

PROBLEM 295. 3 points In the same way, check that the decomposition  $\begin{bmatrix} 3 \\ 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 3 \\ 4 \end{bmatrix}$  is  $\mathbf{y} = \hat{\mathbf{y}} + \boldsymbol{\varepsilon}$  in the regression of  $\mathbf{y} = \begin{bmatrix} 3 \\ 3 \\ 4 \end{bmatrix}$  on  $\mathbf{x}_1 = \begin{bmatrix} 5 \\ 0 \\ 0 \end{bmatrix}$  and  $\mathbf{x}_2 = \begin{bmatrix} -1 \\ 4 \\ 0 \end{bmatrix}$ .

ANSWER. Besides the equation  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\hat{\mathbf{e}}}$  we have to check two things: (1)  $\hat{\mathbf{y}}$  is a linear combination of all the explanatory variables. Since both  $\mathbf{x}_1$  and  $\mathbf{x}_2$  have zero as third coordinate, and they are linearly independent, they span the whole plane, therefore  $\hat{\mathbf{y}}$ , which also has the third coordinate zero, is their linear combination. (2)  $\hat{\hat{\mathbf{e}}}$  is orthogonal to both explanatory variables because its only nonzero coordinate is the third.

The residuals  $\hat{\hat{\mathbf{e}}}$  in the regression on  $\mathbf{x}_1$  are  $\mathbf{y} - \hat{\mathbf{y}} = \begin{bmatrix} 3 \\ 3 \\ 4 \end{bmatrix} - \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \\ 4 \end{bmatrix}$ . This vector is clearly orthogonal to  $\mathbf{x}_1 = \begin{bmatrix} 5 \\ 0 \\ 0 \end{bmatrix}$ . Now let us regress  $\hat{\hat{\mathbf{e}}} = \begin{bmatrix} 0 \\ 3 \\ 4 \end{bmatrix}$  on  $\mathbf{x}_2 = \begin{bmatrix} -1 \\ 4 \\ 0 \end{bmatrix}$ . Say  $\mathbf{h}$  is the vector of fitted values and  $\mathbf{k}$  the residual vector in this regression. We saw in problem 293 that this is the next step in backfitting, but  $\mathbf{k}$  is not the same as the residual vector  $\hat{\mathbf{e}}$  in the multiple regression, because  $\mathbf{k}$  is not orthogonal to  $\mathbf{x}_1$ . In order to get the correct residual in the joint regression and also the correct coefficient of  $\mathbf{x}_2$ , one must regress  $\hat{\hat{\mathbf{e}}}$  only on that part of  $\mathbf{x}_2$  which is orthogonal to  $\mathbf{x}_1$ . This regressor is the dark blue line starting at the origin.

In formulas: One gets the correct  $\hat{\hat{\mathbf{e}}}$  and  $\hat{\beta}_2$  by regressing  $\hat{\hat{\mathbf{e}}} = \mathbf{M}_1\mathbf{y}$  not on  $\mathbf{x}_2$  but on  $\mathbf{M}_1\mathbf{X}_2$ , where  $\mathbf{M}_1 = \mathbf{I} - \mathbf{X}_1(\mathbf{X}_1^\top\mathbf{X}_1)^{-1}\mathbf{X}_1^\top$  is the matrix which forms the residuals under the regression on  $\mathbf{X}_1$ . In other words, one has to remove the influence of  $\mathbf{X}_1$  not only from the dependent but also the independent variables. Instead of regressing the residuals  $\hat{\hat{\mathbf{e}}} = \mathbf{M}_1\mathbf{y}$  on  $\mathbf{X}_2$ , one has to regress them on what is left about  $\mathbf{X}_2$  after we know  $\mathbf{X}_1$ , i.e., on what remains of  $\mathbf{X}_2$  after taking out the effect of  $\mathbf{X}_1$ , which is  $\mathbf{M}_1\mathbf{X}_2$ . The regression which gets the correct  $\hat{\beta}_2$  is therefore

$$(23.0.2) \quad \mathbf{M}_1\mathbf{y} = \mathbf{M}_1\mathbf{X}_2\hat{\beta}_2 + \hat{\hat{\mathbf{e}}}$$

In formulas, the correct  $\hat{\beta}_2$  is

$$(23.0.3) \quad \hat{\beta}_2 = (\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{M}_1\mathbf{y}.$$

This regression also yields the correct covariance matrix. (The only thing which is not right is the number of degrees of freedom). The regression is therefore fully representative of the additional effect of  $\mathbf{x}_2$ , and the plot of  $\hat{\hat{\mathbf{e}}}$  against  $\mathbf{M}_1\mathbf{X}_2$  with the fitted line drawn (which has the correct slope  $\hat{\beta}_2$ ) is called the “added variable plot” for  $\mathbf{X}_2$ . [CW99, pp. 244–246] has a good discussion of added variable plots.

PROBLEM 296. 2 points Show that in the model (23.0.1), the estimator  $\hat{\beta}_2 = (\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{M}_1\mathbf{y}$  is unbiased. Compute  $MSE[\hat{\beta}_2; \beta_2]$ .

ANSWER.  $\hat{\beta}_2 - \beta_2 = (\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{M}_1(\mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \boldsymbol{\varepsilon}) - \beta_2 = (\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{M}_1\boldsymbol{\varepsilon}$  therefore  $MSE[\hat{\beta}_2; \beta_2] = \sigma^2(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{M}_1\mathbf{M}_1^\top\mathbf{X}_2(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1} = \sigma^2(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}$

In order to get an estimate of  $\hat{\beta}_1$ , one can again do what seems intuitive, namely regress  $\mathbf{y} - \mathbf{X}_2\hat{\beta}_2$  on  $\mathbf{X}_1$ . This gives

$$(23.0.4) \quad \hat{\beta}_1 = (\mathbf{X}_1^\top\mathbf{X}_1)^{-1}\mathbf{X}_1^\top(\mathbf{y} - \mathbf{X}_2\hat{\beta}_2).$$

This regression also gives the right residuals, but not the right estimates of the covariance matrix.

PROBLEM 297. *The three Figures in [DM93, p. 22] can be seen in XGobi if you use the instructions in Problem 293. The purple line represents the dependent variable  $\mathbf{y}$ , and the two yellow lines the explanatory variables  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . ( $\mathbf{x}_1$  is the one which is in part red.) The two green lines represent the unconstrained regression  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\boldsymbol{\varepsilon}}$ , and the two red lines the constrained regression  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\boldsymbol{\varepsilon}}$  where  $\mathbf{y}$  is only regressed on  $\mathbf{x}_1$ . The two dark blue lines, barely visible against the dark blue background, represent the regression of  $\mathbf{x}_2$  on  $\mathbf{x}_1$ .*

• a. *The first diagram which XGobi shows on startup is [DM93, diagram (b) on p. 22]. Go into the **Rotation** view and rotate the diagram in such a way that the view is [DM93, Figure (a)]. You may want to delete the two white lines, since they are not shown in Figure (a).*

• b. *Make a geometric argument that the light blue line, which represents  $\hat{\mathbf{y}} - \hat{\mathbf{y}} = \mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})$ , is orthogonal on the green line  $\hat{\boldsymbol{\varepsilon}}$  (this is the green line which ends at the point  $\mathbf{y}$ , i.e., not the green line which starts at the origin).*

ANSWER. The light blue line lies in the plane spanned by  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , and  $\hat{\boldsymbol{\varepsilon}}$  is orthogonal to this plane.  $\square$

• c. *Make a geometric argument that the light blue line is also orthogonal to the red line  $\hat{\mathbf{y}}$  emanating from the origin.*

ANSWER. This is a little trickier. The red line  $\hat{\mathbf{y}}$  is orthogonal to  $\mathbf{x}_1$ , and the green line  $\hat{\boldsymbol{\varepsilon}}$  is also orthogonal to  $\mathbf{x}_1$ . Together,  $\hat{\boldsymbol{\varepsilon}}$  and  $\hat{\mathbf{y}}$  span therefore the plane orthogonal to  $\mathbf{x}_1$ . Since the light blue line lies in the plane spanned by  $\hat{\boldsymbol{\varepsilon}}$  and  $\hat{\mathbf{y}}$ , it is orthogonal to  $\mathbf{x}_1$ .  $\square$

Question 297 shows that the decomposition  $\mathbf{y} = \hat{\mathbf{y}} + (\hat{\mathbf{y}} - \hat{\mathbf{y}}) + \hat{\boldsymbol{\varepsilon}}$  is orthogonal, i.e., all 3 vectors  $\hat{\mathbf{y}}$ ,  $\hat{\mathbf{y}} - \hat{\mathbf{y}}$ , and  $\hat{\boldsymbol{\varepsilon}}$  are orthogonal to each other. This is (22.7.6) in the special case that  $\mathbf{u} = \mathbf{o}$  and therefore  $\boldsymbol{\beta}_0 = \mathbf{o}$ .

One can use this same animation also to show the following: If you first project the purple line on the plane spanned by the yellow lines, you get the green line in the plane. If you then project that green line on  $\mathbf{x}_1$ , which is a subspace of the plane, then you get the red section of the yellow line. This is the same result as if you had projected the purple line directly on  $\mathbf{x}_1$ . A matrix-algebraic proof of this fact is given in (A.6.3).

The same animation allows us to verify the following:

- In the regression of  $\mathbf{y}$  on  $\mathbf{x}_1$ , the coefficient is  $\hat{\beta}_1$ , and the residual is  $\hat{\boldsymbol{\varepsilon}}$ .
- In the regression of  $\mathbf{y}$  on  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , the coefficients are  $\hat{\beta}_1$ ,  $\hat{\beta}_2$ , and the residual is  $\hat{\boldsymbol{\varepsilon}}$ .
- In the regression of  $\mathbf{y}$  on  $\mathbf{x}_1$  and  $\mathbf{M}_1\mathbf{x}_2$ , the coefficients are  $\hat{\beta}_1$ ,  $\hat{\beta}_2$ , and the residual is  $\hat{\boldsymbol{\varepsilon}}$ . The residual is  $\hat{\boldsymbol{\varepsilon}}$  because the space spanned by the regressors

is the same as in the regression on  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , and  $\hat{\boldsymbol{\varepsilon}}$  only depends on the space.

- In the regression of  $\mathbf{y}$  on  $\mathbf{M}_1\mathbf{x}_2$ , the coefficient is  $\hat{\beta}_2$ , because the regressor I am leaving out is orthogonal to  $\mathbf{M}_1\mathbf{x}_2$ . The residual contains the contribution of the left-out variable, i.e., it is  $\hat{\boldsymbol{\varepsilon}} + \hat{\beta}_1\mathbf{x}_1$ .
- But in the regression of  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}_1\mathbf{y}$  on  $\mathbf{M}_1\mathbf{x}_2$ , the coefficient is  $\hat{\beta}_2$  and the residual  $\hat{\boldsymbol{\varepsilon}}$ .

This last statement is (23.0.3).

Now let us turn to proving all this mathematically. The “brute force” proof, i.e., the proof which is conceptually simplest but has to plow through some tedious mathematics, uses (14.2.4) with partitioned matrix inverses. For this we need (23.0.5)

PROBLEM 298. *4 points This is a simplified version of question 393. Show the following, by multiplying  $\mathbf{X}^\top\mathbf{X}$  with its alleged inverse: If  $\mathbf{X} = [\mathbf{X}_1 \quad \mathbf{X}_2]$  has full column rank, then  $(\mathbf{X}^\top\mathbf{X})^{-1}$  is the following partitioned matrix:*

$$(23.0.5) \quad \begin{bmatrix} \mathbf{X}_1^\top\mathbf{X}_1 & \mathbf{X}_1^\top\mathbf{X}_2 \\ \mathbf{X}_2^\top\mathbf{X}_1 & \mathbf{X}_2^\top\mathbf{X}_2 \end{bmatrix}^{-1} = \begin{bmatrix} (\mathbf{X}_1^\top\mathbf{X}_1)^{-1} + \mathbf{K}_1^\top\mathbf{X}_2(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{K}_1 & -\mathbf{K}_1^\top\mathbf{X}_2(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1} \\ -(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{K}_1 & (\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1} \end{bmatrix}$$

where  $\mathbf{M}_1 = \mathbf{I} - \mathbf{X}_1(\mathbf{X}_1^\top\mathbf{X}_1)^{-1}\mathbf{X}_1^\top$  and  $\mathbf{K}_1 = \mathbf{X}_1(\mathbf{X}_1^\top\mathbf{X}_1)^{-1}$ .

From (23.0.5) one sees that the covariance matrix in regression (23.0.3) is the lower left partition of the covariance matrix in the full regression (23.0.1).

PROBLEM 299. *6 points Use the usual formula  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}$  together with (23.0.5) to prove (23.0.3) and (23.0.4).*

ANSWER. (14.2.4) reads here

$$(23.0.6) \quad \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} (\mathbf{X}_1^\top\mathbf{X}_1)^{-1} + \mathbf{K}_1^\top\mathbf{X}_2(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{K}_1 & -\mathbf{K}_1^\top\mathbf{X}_2(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1} \\ -(\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1}\mathbf{X}_2^\top\mathbf{K}_1 & (\mathbf{X}_2^\top\mathbf{M}_1\mathbf{X}_2)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{X}_1^\top\mathbf{y} \\ \mathbf{X}_2^\top\mathbf{y} \end{bmatrix}$$

Since  $M_1 = I - K_1 X_1^\top$ , one can simplify

$$(23.0.7) \quad \hat{\beta}_2 = -(X_2^\top M_1 X_2)^{-1} X_2^\top K_1 X_1^\top \mathbf{y} + (X_2^\top M_1 X_2)^{-1} X_2^\top \mathbf{y}$$

$$(23.0.8) \quad = (X_2^\top M_1 X_2)^{-1} X_2^\top M_1 \mathbf{y}$$

$$(23.0.9) \quad \hat{\beta}_1 = (X_1^\top X_1)^{-1} X_1^\top \mathbf{y} + K_1^\top X_2 (X_2^\top M_1 X_2)^{-1} X_2^\top K_1 X_1^\top \mathbf{y} - K_1^\top X_2 (X_2^\top M_1 X_2)^{-1} X_2^\top \mathbf{y}$$

$$(23.0.10) \quad = K_1^\top \mathbf{y} - K_1^\top X_2 (X_2^\top M_1 X_2)^{-1} X_2^\top (I - K_1 X_1^\top) \mathbf{y}$$

$$(23.0.11) \quad = K_1^\top \mathbf{y} - K_1^\top X_2 (X_2^\top M_1 X_2)^{-1} X_2^\top M_1 \mathbf{y}$$

$$(23.0.12) \quad = K_1^\top (\mathbf{y} - X_2 \hat{\beta}_2)$$

□

[Gre97, pp. 245–7] follow a different proof strategy: he solves the partitioned normal equations

$$(23.0.13) \quad \begin{bmatrix} X_1^\top X_1 & X_1^\top X_2 \\ X_2^\top X_1 & X_2^\top X_2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} X_1^\top \mathbf{y} \\ X_2^\top \mathbf{y} \end{bmatrix}$$

directly, without going through the inverse. A third proof strategy, used by [Seb77, pp. 65–72], is followed in Problems 301 and 302.

PROBLEM 300. 5 points [Gre97, problem 18 on p. 326]. The following matrix gives the slope in the simple regression of the column variable on the row variable:

$$(23.0.14) \quad \begin{array}{cccc} \mathbf{y} & \mathbf{x}_1 & \mathbf{x}_2 & \\ 1 & 0.03 & 0.36 & \mathbf{y} \\ 0.4 & 1 & 0.3 & \mathbf{x}_1 \\ 1.2 & 0.075 & 1 & \mathbf{x}_2 \end{array}$$

For example, if  $\mathbf{y}$  is regressed on  $\mathbf{x}_1$ , the slope is 0.4, but if  $\mathbf{x}_1$  is regressed on  $\mathbf{y}$ , the slope is 0.03. All variables have zero means, so the constant terms in all regressions are zero. What are the two slope coefficients in the multiple regression of  $\mathbf{y}$  on  $\mathbf{x}_1$  and  $\mathbf{x}_2$ ? Hint: Use the partitioned normal equation as given in [Gre97, p. 245] in the special case when each of the partitions of  $\mathbf{X}$  has only one column.

ANSWER.

$$(23.0.15) \quad \begin{bmatrix} x_1^\top x_1 & x_1^\top x_2 \\ x_2^\top x_1 & x_2^\top x_2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} x_1^\top \mathbf{y} \\ x_2^\top \mathbf{y} \end{bmatrix}$$

The first row reads

$$(23.0.16) \quad \hat{\beta}_1 + (x_1^\top x_1)^{-1} x_1^\top x_2 \hat{\beta}_2 = (x_1^\top x_1)^{-1} x_1^\top \mathbf{y}$$

which is the upper line of [Gre97, (6.24) on p. 245], and in our numbers this is  $\hat{\beta}_1 = 0.4 - 0.3$ . The second row reads

$$(23.0.17) \quad (x_2^\top x_2)^{-1} x_2^\top x_1 \hat{\beta}_1 + \hat{\beta}_2 = (x_2^\top x_2)^{-1} x_2^\top \mathbf{y}$$

or in our numbers  $0.075 \hat{\beta}_2 + \hat{\beta}_2 = 1.2$ . Plugging in the formula for  $\hat{\beta}_1$  gives  $0.075 \cdot 0.4 - 0.3 \hat{\beta}_2 + \hat{\beta}_2 = 1.2$ . This gives  $\hat{\beta}_2 = 1.17/0.9775 = 1.196931 = 1.2$  roughly, and  $\hat{\beta}_1 = 0.4 - 0.3 \cdot 0.0409207 = 0.041$  roughly.

PROBLEM 301. Derive (23.0.3) and (23.0.4) from the first order conditions minimizing

$$(23.0.18) \quad (\mathbf{y} - X_1 \beta_1 - X_2 \beta_2)^\top (\mathbf{y} - X_1 \beta_1 - X_2 \beta_2).$$

ANSWER. Start by writing down the OLS objective function for the full model. Perhaps can use the more sophisticated matrix differentiation rules?

$$(23.0.19) \quad (\mathbf{y} - X_1 \beta_1 - X_2 \beta_2)^\top (\mathbf{y} - X_1 \beta_1 - X_2 \beta_2) = \mathbf{y}^\top \mathbf{y} + \beta_1^\top X_1^\top X_1 \beta_1 + \beta_2^\top X_2^\top X_2 \beta_2 - 2\mathbf{y}^\top X_1 \beta_1 - 2\mathbf{y}^\top X_2 \beta_2$$

Taking partial derivatives with respect to  $\beta_1^\top$  and  $\beta_2^\top$  gives

$$(23.0.20) \quad 2\beta_1^\top X_1^\top X_1 - 2\mathbf{y}^\top X_1 + 2\beta_2^\top X_2^\top X_1 \quad \text{or, transposed} \quad 2X_1^\top X_1 \beta_1 - 2X_1^\top \mathbf{y} + 2X_1^\top X_2 \beta_2$$

$$(23.0.21) \quad 2\beta_2^\top X_2^\top X_2 - 2\mathbf{y}^\top X_2 + 2\beta_1^\top X_1^\top X_2 \quad \text{or, transposed} \quad 2X_2^\top X_2 \beta_2 - 2X_2^\top \mathbf{y} + 2X_2^\top X_1 \beta_1$$

Setting them zero and replacing  $\beta_1$  by  $\hat{\beta}_1$  and  $\beta_2$  by  $\hat{\beta}_2$  gives

$$(23.0.22) \quad X_1^\top X_1 \hat{\beta}_1 = X_1^\top (\mathbf{y} - X_2 \hat{\beta}_2)$$

$$(23.0.23) \quad X_2^\top X_2 \hat{\beta}_2 = X_2^\top (\mathbf{y} - X_1 \hat{\beta}_1).$$

Premultiply (23.0.22) by  $X_1(X_1^\top X_1)^{-1}$ :

$$(23.0.24) \quad X_1 \hat{\beta}_1 = X_1 (X_1^\top X_1)^{-1} X_1^\top (\mathbf{y} - X_2 \hat{\beta}_2).$$

Plug this into (23.0.23):

$$(23.0.25) \quad X_2^\top X_2 \hat{\beta}_2 = X_2^\top (\mathbf{y} - X_1 (X_1^\top X_1)^{-1} X_1^\top \mathbf{y} + X_1 (X_1^\top X_1)^{-1} X_1^\top X_2 \hat{\beta}_2)$$

$$(23.0.26) \quad X_2^\top M_1 X_2 \hat{\beta}_2 = X_2^\top M_1 \mathbf{y}.$$

(23.0.26) is the normal equation of the regression of  $M_1 \mathbf{y}$  on  $M_1 X_2$ ; it immediately implies (23.0.1). Once  $\hat{\beta}_2$  is known, (23.0.22) is the normal equation of the regression of  $\mathbf{y} - X_2 \hat{\beta}_2$  on  $X_1$ , which gives (23.0.4).

PROBLEM 302. Using (23.0.3) and (23.0.4) show that the residuals in regressions (23.0.1) are identical to those in the regression of  $M_1 \mathbf{y}$  on  $M_1 X_2$ .

ANSWER.

$$(23.0.27) \quad \hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2$$

$$(23.0.28) \quad = \mathbf{y} - \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top (\mathbf{y} - \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2) - \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2$$

$$(23.0.29) \quad = \mathbf{M}_1 \mathbf{y} - \mathbf{M}_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2. \quad \square$$

PROBLEM 303. The following problem derives one of the main formulas for adding regressors, following [DM93, pp. 19–24]. We are working in model (23.0.1).

• a. 1 point Show that, if  $\mathbf{X}$  has full column rank, then  $\mathbf{X}^\top \mathbf{X}$ ,  $\mathbf{X}_1^\top \mathbf{X}_1$ , and  $\mathbf{X}_2^\top \mathbf{X}_2$  are nonsingular. Hint: A matrix  $\mathbf{X}$  has full column rank if  $\mathbf{X}\mathbf{a} = \mathbf{o}$  implies  $\mathbf{a} = \mathbf{o}$ .

ANSWER. From  $\mathbf{X}^\top \mathbf{X}\mathbf{a} = \mathbf{o}$  follows  $\mathbf{a}^\top \mathbf{X}^\top \mathbf{X}\mathbf{a} = 0$  which can also be written  $\|\mathbf{X}\mathbf{a}\|^2 = 0$ . Therefore  $\mathbf{X}\mathbf{a} = \mathbf{o}$ , and since the columns are linearly independent, it follows  $\mathbf{a} = \mathbf{o}$ .  $\mathbf{X}_1^\top \mathbf{X}_1$  and  $\mathbf{X}_2^\top \mathbf{X}_2$  are nonsingular because, along with  $\mathbf{X}$ , also  $\mathbf{X}_1$  and  $\mathbf{X}_2$  have full column rank.  $\square$

• b. 1 point Define  $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  and  $\mathbf{M}_1 = \mathbf{I} - \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$ . Show that both  $\mathbf{M}$  and  $\mathbf{M}_1$  are projection matrices. (Give the definition of a projection matrix.) Which spaces do they project on? Which space is bigger?

ANSWER. A projection matrix is symmetric and idempotent. That  $\mathbf{M}\mathbf{M} = \mathbf{M}$  is easily verified.  $\mathbf{M}$  projects on the orthogonal complement of the column space of  $\mathbf{X}$ , and  $\mathbf{M}_1$  on that of  $\mathbf{X}_1$ . I.e.,  $\mathbf{M}_1$  projects on the larger space.  $\square$

• c. 2 points Prove that  $\mathbf{M}_1 \mathbf{M} = \mathbf{M}$  and that  $\mathbf{M}\mathbf{X}_1 = \mathbf{O}$  as well as  $\mathbf{M}\mathbf{X}_2 = \mathbf{O}$ . You will need each these equationse below. What is their geometric meaning?

ANSWER.  $\mathbf{X}_1 = \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \mathbf{I} \\ \mathbf{O} \end{bmatrix} = \mathbf{X}\mathbf{A}$ , say. Therefore  $\mathbf{M}_1 \mathbf{M} = (\mathbf{I} - \mathbf{X}\mathbf{A}(\mathbf{A}^\top \mathbf{X}^\top \mathbf{X}\mathbf{A})^{-1} \mathbf{A}^\top \mathbf{X}^\top)$ .

$\mathbf{M}$  because  $\mathbf{X}^\top \mathbf{M} = \mathbf{O}$ . Geometrically this means that the space on which  $\mathbf{M}$  projects is a subspace of the space on which  $\mathbf{M}_1$  projects. To show that  $\mathbf{M}\mathbf{X}_2 = \mathbf{O}$  note that  $\mathbf{X}_2$  can be written in the form  $\mathbf{X}_2 = \mathbf{X}\mathbf{B}$ , too; this time,  $\mathbf{B} = \begin{bmatrix} \mathbf{O} \\ \mathbf{I} \end{bmatrix}$ .  $\mathbf{M}\mathbf{X}_2 = \mathbf{O}$  means geometrically that  $\mathbf{M}$  projects on a space that is orthogonal to all columns of  $\mathbf{X}_2$ .  $\square$

• d. 2 points Show that  $\mathbf{M}_1 \mathbf{X}_2$  has full column rank.

ANSWER. If  $\mathbf{M}_1 \mathbf{X}_2 \mathbf{b} = \mathbf{o}$ , then  $\mathbf{X}_2 \mathbf{b} = \mathbf{X}_1 \mathbf{a}$  for some  $\mathbf{a}$ . We showed this in Problem 196. Therefore  $\begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} -\mathbf{a} \\ \mathbf{b} \end{bmatrix} = \begin{bmatrix} \mathbf{o} \\ \mathbf{o} \end{bmatrix}$ , and since  $\begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix}$  has full column rank, it follows  $\begin{bmatrix} -\mathbf{a} \\ \mathbf{b} \end{bmatrix} = \begin{bmatrix} \mathbf{o} \\ \mathbf{o} \end{bmatrix}$ , in particular  $\mathbf{b} = \mathbf{o}$ .  $\square$

• e. 1 point Here is some more notation: the regression of  $\mathbf{y}$  on  $\mathbf{X}_1$  and  $\mathbf{X}_2$  can also be represented by the equation

$$(23.0.30) \quad \mathbf{y} = \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$$

The difference between (23.0.1) and (23.0.30) is that (23.0.30) contains the parameter estimates, not their true values, and the residuals, not the true disturbances. Explain the difference between residuals and disturbances, and between the fitted regression line and the true regression line.

• f. 1 point Verify that premultiplication of (23.0.30) by  $\mathbf{M}_1$  gives

$$(23.0.31) \quad \mathbf{M}_1 \mathbf{y} = \mathbf{M}_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$$

ANSWER. We need  $\mathbf{M}_1 \mathbf{X}_1 = \mathbf{O}$  and  $\mathbf{M}_1 \hat{\boldsymbol{\varepsilon}} = \mathbf{M}_1 \mathbf{M}\mathbf{y} = \mathbf{M}\mathbf{y} = \hat{\boldsymbol{\varepsilon}}$  (or this can also be seen because  $\mathbf{X}_1^\top \hat{\boldsymbol{\varepsilon}} = \mathbf{o}$ ).

• g. 2 points Prove that (23.0.31) is the fit which one gets if one regresses  $\mathbf{M}_1 \mathbf{y}$  on  $\mathbf{M}_1 \mathbf{X}_2$ . In other words, if one runs OLS with dependent variable  $\mathbf{M}_1 \mathbf{y}$  and explanatory variables  $\mathbf{M}_1 \mathbf{X}_2$ , one gets the same  $\hat{\boldsymbol{\beta}}_2$  and  $\hat{\boldsymbol{\varepsilon}}$  as in (23.0.31), which are the same  $\hat{\boldsymbol{\beta}}_2$  and  $\hat{\boldsymbol{\varepsilon}}$  as in the complete regression (23.0.30).

ANSWER. According to Problem ?? we have to check  $\mathbf{X}_2^\top \mathbf{M}_1 \hat{\boldsymbol{\varepsilon}} = \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{M}\mathbf{y} = \mathbf{X}_2^\top \mathbf{M}\mathbf{y}$  and  $\mathbf{O}\mathbf{y} = \mathbf{o}$ .

• h. 1 point Show that  $\mathcal{V}[\hat{\boldsymbol{\beta}}_2] = (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1}$ . Are the variance estimates and confidence intervals valid, which the computer automatically prints out if one regresses  $\mathbf{M}_1 \mathbf{y}$  on  $\mathbf{M}_1 \mathbf{X}_2$ ?

ANSWER. Yes except for the number of degrees of freedom.

• i. 4 points If one premultiplies (23.0.1) by  $\mathbf{M}_1$ , one obtains

$$(23.0.32) \quad \mathbf{M}_1 \mathbf{y} = \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{M}_1 \boldsymbol{\varepsilon}, \quad \mathbf{M}_1 \boldsymbol{\varepsilon} \sim (\mathbf{o}, \sigma^2 \mathbf{M}_1)$$

Although the covariance matrix of the disturbance  $\mathbf{M}_1 \boldsymbol{\varepsilon}$  in (23.0.32) is no longer spherical, show that nevertheless the  $\hat{\boldsymbol{\beta}}_2$  obtained by running OLS on (23.0.32) is BLUE of  $\boldsymbol{\beta}_2$  based on the information given in (23.0.32) (i.e., assuming that  $\mathbf{M}_1 \mathbf{y}$  and  $\mathbf{M}_1 \mathbf{X}_2$  are known, but not necessarily  $\mathbf{M}_1$ ,  $\mathbf{y}$ , and  $\mathbf{X}_2$  separately). Hint: the proof is almost identical to the proof that for spherically distributed disturbances OLS is BLUE (e.g. given in [DM93, p. 159]), but you have to add some  $\mathbf{M}_1$ 's to your formulas.

ANSWER. Any other linear estimator  $\tilde{\boldsymbol{\gamma}}$  of  $\boldsymbol{\beta}_2$  can be written as  $\tilde{\boldsymbol{\gamma}} = ((\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{C}) \mathbf{M}_1 \mathbf{y}$ . Its expected value is  $E[\tilde{\boldsymbol{\gamma}}] = (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{C} \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{C} \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2$ . For  $\tilde{\boldsymbol{\gamma}}$  to be unbiased, regardless of the value of  $\boldsymbol{\beta}_2$ ,  $\mathbf{C}$  must satisfy  $\mathbf{C} \mathbf{M}_1 \mathbf{X}_2 = \mathbf{O}$ . From this follows  $\mathcal{MSE}[\tilde{\boldsymbol{\gamma}}; \boldsymbol{\beta}_2] = \sigma^2 ((\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{C}) \mathbf{M}_1 (\mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} + \mathbf{C}^\top) = \sigma^2 (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} + \sigma^2 \mathbf{C} \mathbf{M}_1 \mathbf{C}^\top$ . i.e., it exceeds the  $\mathcal{MSE}$ -matrix of  $\hat{\boldsymbol{\beta}}$  by a nonnegative definite matrix. Is it unique? The form



for the BLUE is not unique, since one can add any  $C$  with  $CM_1C^T = O$  or equivalently  $CM_1 = O$  or  $C = AX$  for some  $A$ . However such a  $C$  applied to a dependent variable of the form  $M_1\mathbf{y}$  will give the null vector, therefore the *values* of the BLUE for those values of  $\mathbf{y}$  which are possible are indeed unique.  $\square$

• j. 1 point Once  $\hat{\beta}_2$  is known, one can move it to the left hand side in (23.0.30) to get

$$(23.0.33) \quad \mathbf{y} - \mathbf{X}_2\hat{\beta}_2 = \mathbf{X}_1\hat{\beta}_1 + \hat{\varepsilon}$$

Prove that one gets the right values of  $\hat{\beta}_1$  and of  $\hat{\varepsilon}$  if one regresses  $\mathbf{y} - \mathbf{X}_2\hat{\beta}_2$  on  $\mathbf{X}_1$ .

ANSWER. The simplest answer just observes that  $\mathbf{X}_1^T\hat{\varepsilon} = \mathbf{o}$ . Or: The normal equation for this pseudo-regression is  $\mathbf{X}_1^T\mathbf{y} - \mathbf{X}_1^T\mathbf{X}_2\hat{\beta}_2 = \mathbf{X}_1^T\mathbf{X}_1\hat{\beta}_1$ , which holds due to the normal equation for the full model.  $\square$

• k. 1 point Does (23.0.33) also give the right covariance matrix for  $\hat{\beta}_1$ ?

ANSWER. No, since  $\mathbf{y} - \mathbf{X}_2\hat{\beta}_2$  has a different covariance matrix than  $\sigma^2\mathbf{I}$ .  $\square$

The following Problems gives some applications of the results in Problem 303. You are allowed to use the results of Problem 303 without proof.

PROBLEM 304. Assume your regression involves an intercept, i.e., the matrix of regressors is  $[\boldsymbol{\iota} \quad \mathbf{X}]$ , where  $\mathbf{X}$  is the matrix of the “true” explanatory variables with no vector of ones built in, and  $\boldsymbol{\iota}$  the vector of ones. The regression can therefore be written

$$(23.0.34) \quad \mathbf{y} = \boldsymbol{\iota}\alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

• a. 1 point Show that the OLS estimate of the slope parameters  $\boldsymbol{\beta}$  can be obtained by regressing  $\underline{\mathbf{y}}$  on  $\underline{\mathbf{X}}$  without intercept, where  $\underline{\mathbf{y}}$  and  $\underline{\mathbf{X}}$  are the variables with their means taken out, i.e.,  $\underline{\mathbf{y}} = \mathbf{D}\mathbf{y}$  and  $\underline{\mathbf{X}} = \mathbf{D}\mathbf{X}$ , with  $\mathbf{D} = \mathbf{I} - \frac{1}{n}\boldsymbol{\iota}\boldsymbol{\iota}^T$ .

ANSWER. This is called the “sweeping out of means.” It follows immediately from (23.0.3). This is the usual procedure to do regression with a constant term: in simple regression  $y_i = \alpha + \beta x_i + \varepsilon_i$ , (23.0.3) is equation (14.2.22):

$$(23.0.35) \quad \hat{\beta} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}.$$

• b. Show that the OLS estimate of the intercept is  $\hat{\alpha} = \bar{y} - \bar{\mathbf{x}}^T\hat{\boldsymbol{\beta}}$  where  $\bar{\mathbf{x}}^T$  is the row vector of column means of  $\mathbf{X}$ , i.e.,  $\bar{\mathbf{x}}^T = \frac{1}{n}\boldsymbol{\iota}^T\mathbf{X}$ .  $\square$

ANSWER. This is exactly (23.0.4). Here is a more specific argument: The intercept  $\hat{\alpha}$  is obtained by regressing  $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  on  $\boldsymbol{\iota}$ . The normal equation for this second regression is  $\boldsymbol{\iota}^T\mathbf{y} - \boldsymbol{\iota}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \boldsymbol{\iota}^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ . If  $\bar{y}$  is the mean of  $\mathbf{y}$ , and  $\bar{\mathbf{x}}^T$  the row vector consisting of means of the columns of  $\mathbf{X}$ , then this gives  $\bar{y} = \bar{\mathbf{x}}^T\hat{\boldsymbol{\beta}} + \hat{\alpha}$ . In the case of simple regression, this was derived earlier as formula (14.2.23).

• c. 2 points Show that  $MSE[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}] = \sigma^2(\underline{\mathbf{X}}^T\underline{\mathbf{X}})^{-1}$ . (Use the formula for  $\hat{\boldsymbol{\beta}}$ .)

ANSWER. Since

$$(23.0.36) \quad \begin{bmatrix} \boldsymbol{\iota}^T \\ \mathbf{X}^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\iota} & \mathbf{X} \end{bmatrix} = \begin{bmatrix} n & n\bar{\mathbf{x}}^T \\ \bar{\mathbf{x}}n & \mathbf{X}^T\mathbf{X} \end{bmatrix},$$

it follows by Problem 393

$$(23.0.37) \quad \left( \begin{bmatrix} \boldsymbol{\iota}^T \\ \mathbf{X}^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\iota} & \mathbf{X} \end{bmatrix} \right)^{-1} = \begin{bmatrix} 1/n + \bar{\mathbf{x}}^T(\mathbf{X}^T\mathbf{X})^{-1}\bar{\mathbf{x}} & -\bar{\mathbf{x}}^T(\mathbf{X}^T\mathbf{X})^{-1} \\ -(\mathbf{X}^T\mathbf{X})^{-1}\bar{\mathbf{x}} & (\mathbf{X}^T\mathbf{X})^{-1} \end{bmatrix}$$

In other words, one simply does as if the actual regressors had been the data with their means removed, and then takes the inverse of that design matrix. The only place where one has to be careful is the number of degrees of freedom. See also Seber [Seb77, section 11.7] about centering and scaling the data.

• d. 3 points Show that  $\hat{\mathbf{y}} - \boldsymbol{\iota}\bar{y} = \underline{\mathbf{X}}\hat{\boldsymbol{\beta}}$ .

ANSWER. First note that  $\mathbf{X} = \underline{\mathbf{X}} + \frac{1}{n}\boldsymbol{\iota}\boldsymbol{\iota}^T\mathbf{X} = \underline{\mathbf{X}} + \boldsymbol{\iota}\bar{\mathbf{x}}^T$  where  $\bar{\mathbf{x}}^T$  is the row vector of means of  $\mathbf{X}$ . By definition,  $\hat{\mathbf{y}} = \boldsymbol{\iota}\hat{\alpha} + \mathbf{X}\hat{\boldsymbol{\beta}} = \boldsymbol{\iota}\hat{\alpha} + \underline{\mathbf{X}}\hat{\boldsymbol{\beta}} + \boldsymbol{\iota}\bar{\mathbf{x}}^T\hat{\boldsymbol{\beta}} = \boldsymbol{\iota}(\hat{\alpha} + \bar{\mathbf{x}}^T\hat{\boldsymbol{\beta}}) + \underline{\mathbf{X}}\hat{\boldsymbol{\beta}} = \boldsymbol{\iota}\bar{y} + \underline{\mathbf{X}}\hat{\boldsymbol{\beta}}$ .

• e. 2 points Show that  $R^2 = \frac{\underline{\mathbf{y}}^T\underline{\mathbf{X}}(\underline{\mathbf{X}}^T\underline{\mathbf{X}})^{-1}\underline{\mathbf{X}}^T\underline{\mathbf{y}}}{\underline{\mathbf{y}}^T\underline{\mathbf{y}}}$

ANSWER.

$$(23.0.38) \quad R^2 = \frac{(\hat{\mathbf{y}} - \bar{y}\boldsymbol{\iota})^T(\hat{\mathbf{y}} - \bar{y}\boldsymbol{\iota})}{\underline{\mathbf{y}}^T\underline{\mathbf{y}}} = \frac{\hat{\boldsymbol{\beta}}^T\underline{\mathbf{X}}^T\underline{\mathbf{X}}\hat{\boldsymbol{\beta}}}{\underline{\mathbf{y}}^T\underline{\mathbf{y}}}$$

and now plugging in the formula for  $\hat{\boldsymbol{\beta}}$  the result follows.

• f. 3 points Now, split once more  $\underline{\mathbf{X}} = [\underline{\mathbf{X}}_1 \quad \underline{\mathbf{x}}_2]$  where the second part  $\underline{\mathbf{x}}_2$  consists of one column only, and  $\underline{\mathbf{X}}_1$  is, as above, the  $\mathbf{X}$  matrix with the column means taken out. Conformably,  $\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}$ . Show that

$$(23.0.39) \quad \text{var}[\hat{\beta}_2] = \frac{\sigma^2}{\underline{\mathbf{x}}_2^T\underline{\mathbf{x}}_2(1 - R_2^2)}$$

where  $R_2^2$  is the  $R^2$  in the regression of  $\underline{\mathbf{x}}_2$  on all other variables in  $\underline{\mathbf{X}}$ . This is [Gre97, (9.3) on p. 421]. Hint: you should first show that  $\text{var}[\hat{\beta}_2] = \sigma^2/\underline{\mathbf{x}}_2^T\underline{\mathbf{M}}_1\underline{\mathbf{x}}_2$  where  $\underline{\mathbf{M}}_1 = \mathbf{I} - \underline{\mathbf{X}}_1(\underline{\mathbf{X}}_1^T\underline{\mathbf{X}}_1)^{-1}\underline{\mathbf{X}}_1^T$ . Here is an interpretation of (23.0.39) which you don't have to prove:  $\sigma^2/\underline{\mathbf{x}}_2^T\underline{\mathbf{x}}_2$  is the variance in a simple regression with a constant

term and  $\mathbf{x}_2$  as the only explanatory variable, and  $1/(1 - R_2^2)$  is called the variance inflation factor.

ANSWER. Note that we are not talking about the variance of the constant term but that of all the other terms.

$$(23.0.40) \quad \mathbf{x}_2^\top \underline{\mathbf{M}}_1 \mathbf{x}_2 = \mathbf{x}_2^\top \mathbf{x}_2 + \mathbf{x}_2^\top \underline{\mathbf{X}}_1 (\underline{\mathbf{X}}_1^\top \underline{\mathbf{X}}_1)^{-1} \underline{\mathbf{X}}_1^\top \mathbf{x}_2 = \mathbf{x}_2^\top \mathbf{x}_2 \left( 1 + \frac{\mathbf{x}_2^\top \underline{\mathbf{X}}_1 (\underline{\mathbf{X}}_1^\top \underline{\mathbf{X}}_1)^{-1} \underline{\mathbf{X}}_1^\top \mathbf{x}_2}{\mathbf{x}_2^\top \mathbf{x}_2} \right)$$

and since the fraction is  $R_2^2$ , i.e., it is the  $R^2$  in the regression of  $\mathbf{x}_2$  on all other variables in  $\underline{\mathbf{X}}$ , we get the result.  $\square$

## CHAPTER 24

## Residuals: Standardized, Predictive, “Studentized”

## 24.1. Three Decisions about Plotting Residuals

After running a regression it is always advisable to look at the residuals. Here one has to make three decisions.

The first decision is whether to look at the ordinary residuals

$$(24.1.1) \quad \hat{\varepsilon}_i = \mathbf{y}_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$$

( $\mathbf{x}_i^\top$  is the  $i$ th row of  $\mathbf{X}$ ), or the “predictive” residuals, which are the residuals computed using the OLS estimate of  $\boldsymbol{\beta}$  gained from all the other data except the data point where the residual is taken. If one writes  $\hat{\boldsymbol{\beta}}(i)$  for the OLS estimate without the  $i$ th observation, the defining equation for the  $i$ th predictive residual, which we call  $\hat{\varepsilon}_i(i)$ , is

$$(24.1.2) \quad \hat{\varepsilon}_i(i) = \mathbf{y}_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}(i).$$

The second decision is whether to standardize the residuals or not, i.e., whether to divide them by their estimated standard deviations or not. Since  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}\mathbf{y}$ , the variance of the  $i$ th ordinary residual is

$$(24.1.3) \quad \text{var}[\hat{\varepsilon}_i] = \sigma^2 m_{ii} = \sigma^2(1 - h_{ii}),$$

and regarding the predictive residuals it will be shown below, see (24.2.9), that

$$(24.1.4) \quad \text{var}[\hat{\varepsilon}_i(i)] = \frac{\sigma^2}{m_{ii}} = \frac{\sigma^2}{1 - h_{ii}}.$$

Here

$$(24.1.5) \quad h_{ii} = \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i.$$

(Note that  $\mathbf{x}_i$  is the  $i$ th row of  $\mathbf{X}$  written as a column vector.)  $h_{ii}$  is the  $i$ th diagonal element of the “hat matrix”  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ , the projector on the column space of  $\mathbf{X}$ . This projector is called “hat matrix” because  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ , i.e.,  $\mathbf{H}$  puts the “hat” on  $\mathbf{y}$ .

PROBLEM 305. 2 points Show that the  $i$ th diagonal element of the “hat matrix”  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is  $\mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i$  where  $\mathbf{x}_i$  is the  $i$ th row of  $\mathbf{X}$  written as a column vector.

ANSWER. In terms of  $\mathbf{e}_i$ , the  $n$ -vector with 1 on the  $i$ th place and 0 everywhere else,  $\mathbf{X}^\top \mathbf{e}_i$ , and the  $i$ th diagonal element of the hat matrix is  $\mathbf{e}_i^\top \mathbf{H} \mathbf{e}_i = \mathbf{e}_i^\top \mathbf{X}_i (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_i^\top \mathbf{e}_i = \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i$ .

PROBLEM 306. 2 points The variance of the  $i$ th disturbance is  $\sigma^2$ . Is the variance of the  $i$ th residual bigger than  $\sigma^2$ , smaller than  $\sigma^2$ , or equal to  $\sigma^2$ ? (Before doing math, first argue in words what you would expect it to be.) What about the variance of the predictive residual? Prove your answers mathematically. You are allowed to use (24.2.9) without proof.

ANSWER. Here is only the math part of the answer:  $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}\mathbf{y}$ . Since  $\mathbf{M} = \mathbf{I} - \mathbf{H}$  is idempotent and symmetric, we get  $\mathcal{V}[\mathbf{M}\mathbf{y}] = \sigma^2 \mathbf{M}$ , in particular this means  $\text{var}[\hat{\varepsilon}_i] = \sigma^2 m_{ii}$  where  $m_{ii}$  is the  $i$ th diagonal element of  $\mathbf{M}$ . Then  $m_{ii} = 1 - h_{ii}$ . Since all diagonal elements of projection matrix are between 0 and 1, the answer is: the variances of the ordinary residuals cannot be bigger than  $\sigma^2$ . Regarding predictive residuals, if we plug  $m_{ii} = 1 - h_{ii}$  into (24.2.9) it becomes

$$(24.1.6) \quad \hat{\varepsilon}_i(i) = \frac{1}{m_{ii}} \hat{\varepsilon}_i \quad \text{therefore} \quad \text{var}[\hat{\varepsilon}_i(i)] = \frac{1}{m_{ii}^2} \sigma^2 m_{ii} = \frac{\sigma^2}{m_{ii}}$$

which is bigger than  $\sigma^2$ .

PROBLEM 307. Decide in the following situations whether you want predictive residuals or ordinary residuals, and whether you want them standardized or not.

• a. 1 point You are looking at the residuals in order to check whether the associated data points are outliers and do perhaps not belong into the model.

ANSWER. Here one should use the predictive residuals. If the  $i$ th observation is an outlier which should not be in the regression, then one should not use it when running the regression. Inclusion may have a strong influence on the regression result, and therefore the residual may be as conspicuous. One should standardize them.

• b. 1 point You are looking at the residuals in order to assess whether there is heteroskedasticity.

ANSWER. Here you want them standardized, but there is no reason to use the predictive residuals. Ordinary residuals are a little more precise than predictive residuals because they are based on more observations.

• c. 1 point You are looking at the residuals in order to assess whether the disturbances are autocorrelated.

ANSWER. Same answer as for b.

• d. 1 point You are looking at the residuals in order to assess whether the disturbances are normally distributed.

ANSWER. In my view, one should make a normal QQ-plot of standardized residuals, but one should not use the predictive residuals. To see why, let us first look at the distribution of the standardized residuals before division by  $s$ . Each  $\hat{\varepsilon}_i/\sqrt{1-h_{ii}}$  is normally distributed with mean zero and standard deviation  $\sigma$ . (But different such residuals are not independent.) If one takes a QQ-plot of those residuals against the normal distribution, one will get in the limit a straight line with slope  $\sigma$ . If one divides every residual by  $s$ , the slope will be close to 1, but one will again get something approximating a straight line. The fact that  $s$  is random does not affect the relation of the residuals to each other, and this relation is what determines whether or not the QQ-plot approximates a straight line.

But Belsley, Kuh, and Welsch on [BKW80, p. 43] draw a normal probability plot of the studentized, not the standardized, residuals. They give no justification for their choice. I think it is the wrong choice. □

• e. 1 point Is there any situation in which you do not want to standardize the residuals?

ANSWER. Standardization is a mathematical procedure which is justified when certain conditions hold. But there is no guarantee that these conditions actually hold, and in order to get a more immediate impression of the fit of the curve one may want to look at the unstandardized residuals. □

The third decision is how to plot the residuals. Never do it against  $\mathbf{y}$ . Either do it against the predicted  $\hat{\mathbf{y}}$ , or make several plots against all the columns of the  $\mathbf{X}$ -matrix.

In time series, also a plot of the residuals against time is called for.

Another option are the partial residual plots, see about this also (23.0.2). Say  $\hat{\boldsymbol{\beta}}[h]$  is the estimated parameter vector, which is estimated with the full model, but after estimation we drop the  $h$ -th parameter, and  $\mathbf{X}[h]$  is the  $\mathbf{X}$ -matrix without the  $h$ th column, and  $\mathbf{x}_h$  is the  $h$ th column of the  $\mathbf{X}$ -matrix. Then by (23.0.4), the estimate of the  $h$ th slope parameter is the same as that in the simple regression of  $\mathbf{y} - \mathbf{X}[h]\hat{\boldsymbol{\beta}}[h]$  on  $\mathbf{x}_h$ . The plot of  $\mathbf{y} - \mathbf{X}[h]\hat{\boldsymbol{\beta}}[h]$  against  $\mathbf{x}_h$  is called the  $h$ th partial residual plot.

To understand this better, start out with a regression  $\mathbf{y}_i = \alpha + \beta x_i + \gamma z_i + \varepsilon_i$ ; which gives you the fitted values  $\mathbf{y}_i = \hat{\alpha} + \hat{\beta}x_i + \hat{\gamma}z_i + \hat{\varepsilon}_i$ . Now if you regress  $\mathbf{y}_i - \hat{\alpha} - \hat{\beta}x_i$  on  $x_i$  and  $z_i$  then the intercept will be zero and the estimated coefficient of  $x_i$  will be zero, and the estimated coefficient of  $z_i$  will be  $\hat{\gamma}$ , and the residuals will be  $\hat{\varepsilon}_i$ . The plot of  $\mathbf{y}_i - \hat{\alpha} - \hat{\beta}x_i$  versus  $z_i$  is the partial residuals plot for  $z$ .

### 24.2. Relationship between Ordinary and Predictive Residuals

In equation (24.1.2), the  $i$ th predictive residuals was defined in terms of  $\hat{\boldsymbol{\beta}}(i)$ , the parameter estimate from the regression of  $\mathbf{y}$  on  $\mathbf{X}$  with the  $i$ th observation left out. We will show now that there is a very simple mathematical relationship between the  $i$ th predictive residual and the  $i$ th ordinary residual, namely, equation (24.2.9).

(It is therefore not necessary to run  $n$  different regressions to get the  $n$  predictive residuals.)

We will write  $\mathbf{y}(i)$  for the  $\mathbf{y}$  vector with the  $i$ th element deleted, and  $\mathbf{X}(i)$  is the matrix  $\mathbf{X}$  with the  $i$ th row deleted.

PROBLEM 308. 2 points Show that

$$(24.2.1) \quad \mathbf{X}(i)^\top \mathbf{X}(i) = \mathbf{X}^\top \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^\top$$

$$(24.2.2) \quad \mathbf{X}(i)^\top \mathbf{y}(i) = \mathbf{X}^\top \mathbf{y} - \mathbf{x}_i y_i.$$

ANSWER. Write (24.2.2) as  $\mathbf{X}^\top \mathbf{y} = \mathbf{X}(i)^\top \mathbf{y}(i) + \mathbf{x}_i y_i$ , and observe that with our definition  $\mathbf{x}_i$  as column vectors representing the rows of  $\mathbf{X}$ ,  $\mathbf{X}^\top = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]$ . Therefore

$$(24.2.3) \quad \mathbf{X}^\top \mathbf{y} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n] \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \mathbf{x}_1 y_1 + \dots + \mathbf{x}_n y_n.$$

An important stepping stone towards the proof of (24.2.9) is equation (24.2.4) which gives a relationship between  $h_{ii}$  and

$$(24.2.4) \quad h_{ii}(i) = \mathbf{x}_i^\top (\mathbf{X}(i)^\top \mathbf{X}(i))^{-1} \mathbf{x}_i.$$

$\hat{y}_i(i) = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}(i)$  has variance  $\sigma^2 h_{ii}(i)$ . The following problems give the steps necessary to prove (24.2.8). We begin with a simplified version of theorem A.8.2 in the Mathematical Appendix:

THEOREM 24.2.1. Let  $\mathbf{A}$  be a nonsingular  $k \times k$  matrix,  $\delta \neq 0$  a scalar, and  $\mathbf{b}$  a  $k \times 1$  vector with  $\mathbf{b}^\top \mathbf{A}^{-1} \mathbf{b} + \delta \neq 0$ . Then

$$(24.2.5) \quad \left( \mathbf{A} + \frac{\mathbf{b} \mathbf{b}^\top}{\delta} \right)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{b} \mathbf{b}^\top \mathbf{A}^{-1}}{\delta + \mathbf{b}^\top \mathbf{A}^{-1} \mathbf{b}}.$$

PROBLEM 309. Prove (24.2.5) by showing that the product of the matrix with alleged inverse is the unit matrix.

PROBLEM 310. As an application of (24.2.5) show that

$$(24.2.6) \quad (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1}}{1 - h_{ii}} \quad \text{is the inverse of} \quad \mathbf{X}(i)^\top \mathbf{X}(i).$$

ANSWER. This is (24.2.5), or (A.8.20), with  $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$ ,  $\mathbf{b} = \mathbf{x}_i$ , and  $\delta = -1$ .

PROBLEM 311. Using (24.2.6) show that

$$(24.2.7) \quad (\mathbf{X}(i)^\top \mathbf{X}(i))^{-1} \mathbf{x}_i = \frac{1}{1 - h_{ii}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i,$$

and using (24.2.7) show that  $h_{ii}(i)$  is related to  $h_{ii}$  by the equation

$$(24.2.8) \quad 1 + h_{ii}(i) = \frac{1}{1 - h_{ii}}$$

[Gre97, (9-37) on p. 445] was apparently not aware of this relationship.

PROBLEM 312. Prove the following mathematical relationship between predictive residuals and ordinary residuals:

$$(24.2.9) \quad \hat{\varepsilon}_i(i) = \frac{1}{1 - h_{ii}} \hat{\varepsilon}_i$$

which is the same as (21.0.29), only in a different notation.

ANSWER. For this we have to apply the above mathematical tools. With the help of (24.2.7) (transpose it!) and (24.2.2), (24.1.2) becomes

$$\begin{aligned} \hat{\varepsilon}_i(i) &= \mathbf{y}_i - \mathbf{x}_i^\top (\mathbf{X}(i)^\top \mathbf{X}(i))^{-1} \mathbf{X}(i)^\top \mathbf{y}(i) \\ &= \mathbf{y}_i - \frac{1}{1 - h_{ii}} \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{y} - \mathbf{x}_i \mathbf{y}_i) \\ &= \mathbf{y}_i - \frac{1}{1 - h_{ii}} \mathbf{x}_i^\top \hat{\boldsymbol{\beta}} + \frac{1}{1 - h_{ii}} \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{y}_i \\ &= \mathbf{y}_i \left( 1 + \frac{h_{ii}}{1 - h_{ii}} \right) - \frac{1}{1 - h_{ii}} \mathbf{x}_i^\top \hat{\boldsymbol{\beta}} \\ &= \frac{1}{1 - h_{ii}} (\mathbf{y}_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}) \end{aligned}$$

This is a little tedious but simplifies extremely nicely at the end.  $\square$

The relationship (24.2.9) is so simple because the estimation of  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$  can be done in two steps. First collect the information which the  $n - 1$  observations other than the  $i$ th contribute to the estimation of  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$  is contained in  $\hat{\mathbf{y}}_i(i)$ . The information from all observations except the  $i$ th can be written as

$$(24.2.10) \quad \hat{\mathbf{y}}_i(i) = \eta_i + \boldsymbol{\delta}_i \quad \boldsymbol{\delta}_i \sim (0, \sigma^2 h_{ii}(i))$$

Here  $\boldsymbol{\delta}_i$  is the “sampling error” or “estimation error”  $\hat{\mathbf{y}}_i(i) - \eta_i$  from the regression of  $\mathbf{y}(i)$  on  $\mathbf{X}(i)$ . If we combine this compound “observation” with the  $i$ th observation  $\mathbf{y}_i$ , we get

$$(24.2.11) \quad \begin{bmatrix} \hat{\mathbf{y}}_i(i) \\ \mathbf{y}_i \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \eta_i + \begin{bmatrix} \boldsymbol{\delta}_i \\ \boldsymbol{\varepsilon}_i \end{bmatrix} \quad \begin{bmatrix} \boldsymbol{\delta}_i \\ \boldsymbol{\varepsilon}_i \end{bmatrix} \sim \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma^2 \begin{bmatrix} h_{ii}(i) & 0 \\ 0 & 1 \end{bmatrix} \right)$$

This is a regression model similar to model (14.1.1), but this time with a nonspherical covariance matrix.

PROBLEM 313. Show that the BLUE of  $\eta_i$  in model (24.2.11) is

$$(24.2.12) \quad \hat{\mathbf{y}}_i = (1 - h_{ii}) \hat{\mathbf{y}}_i(i) + h_{ii} \mathbf{y}_i = \hat{\mathbf{y}}_i(i) + h_{ii} \hat{\varepsilon}_i(i)$$

Hint: apply (24.2.8). Use this to prove (24.2.9).

ANSWER. As shown in problem 178, the BLUE in this situation is the weighted average of observations with the weights proportional to the inverses of the variances. I.e., the first observation has weight

$$(24.2.13) \quad \frac{1/h_{ii}(i)}{1/h_{ii}(i) + 1} = \frac{1}{1 + h_{ii}(i)} = 1 - h_{ii}.$$

Since the sum of the weights must be 1, the weight of the second observation is  $h_{ii}$ .

Here is an alternative solution, using formula (19.0.6) for the BLUE, which reads here

$$\begin{aligned} \hat{\mathbf{y}}_i &= \left( \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{h_{ii}}{1-h_{ii}} & 0 \\ 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{h_{ii}}{1-h_{ii}} & 0 \\ 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} \hat{\mathbf{y}}_i(i) \\ \mathbf{y}_i \end{bmatrix} = \\ &= h_{ii} \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{1-h_{ii}}{h_{ii}} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \hat{\mathbf{y}}_i(i) \\ \mathbf{y}_i \end{bmatrix} = (1 - h_{ii}) \hat{\mathbf{y}}_i(i) + h_{ii} \mathbf{y}_i. \end{aligned}$$

Now subtract this last formula from  $\mathbf{y}_i$  to get  $\mathbf{y}_i - \hat{\mathbf{y}}_i = (1 - h_{ii})(\mathbf{y}_i - \hat{\mathbf{y}}_i(i))$ , which is (24.2.9).

### 24.3. Standardization

In this section we will show that the standardized predictive residual is what sometimes called the “studentized” residual. It is recommended not to use the term “studentized residual” but say “standardized predictive residual” instead.

The standardization of the ordinary residuals has two steps: every  $\hat{\varepsilon}_i$  is divided by its “relative” standard deviation  $\sqrt{1 - h_{ii}}$ , and then by  $s$ , an estimate of  $\sigma$ , the standard deviation of the true disturbances. In formulas,

$$(24.3.1) \quad \text{the } i\text{th standardized ordinary residual} = \frac{\hat{\varepsilon}_i}{s\sqrt{1 - h_{ii}}}.$$

Standardization of the  $i$ th predictive residual has the same two steps: first divide the predictive residual (24.2.9) by the relative standard deviation, and then divide by  $s(i)$ . But a look at formula (24.2.9) shows that the ordinary and the predictive residuals differ only by a nonrandom factor. Therefore the first step of the standardization yields exactly the same result whether one starts with an ordinary or a predictive residual. Standardized predictive residuals differ therefore from standardized ordinary residuals only in the second step:

$$(24.3.2) \quad \text{the } i\text{th standardized predictive residual} = \frac{\hat{\varepsilon}_i}{s(i)\sqrt{1 - h_{ii}}}.$$

Note that equation (24.3.2) writes the standardized predictive residual as a function of the ordinary residual, not the predictive residual. The standardized predictive residual is sometimes called the “studentized” residual.

PROBLEM 314. 3 points The  $i$ th predictive residual has the formula

$$(24.3.3) \quad \hat{\varepsilon}_i(i) = \frac{1}{1 - h_{ii}} \hat{\varepsilon}_i$$

You do not have to prove this formula, but you are asked to derive the standard deviation of  $\hat{\varepsilon}_i(i)$ , and to derive from it a formula for the standardized  $i$ th predictive residual.

This similarity between these two formulas has led to widespread confusion. Even [BKW80] seem to have been unaware of the significance of “studentization”; they do not work with the concept of predictive residuals at all.

The standardized predictive residuals have a  $t$ -distribution, because they are a normally distributed variable divided by an independent  $\chi^2$  over its degrees of freedom. (But note that the joint distribution of all standardized predictive residuals is *not* a multivariate  $t$ .) Therefore one can use the quantiles of the  $t$ -distribution to judge, from the size of these residuals, whether one has an extreme observation or not.

PROBLEM 315. Following [DM93, p. 34], we will use (23.0.3) and the other formulas regarding additional regressors to prove the following: If you add a dummy variable which has the value 1 for the  $i$ th observation and the value 0 for all other observations to your regression, then the coefficient estimate of this dummy is the  $i$ th predictive residual, and the coefficient estimate of the other parameters after inclusion of this dummy is equal to  $\hat{\beta}(i)$ . To fix notation (and without loss of generality), assume the  $i$ th observation is the last observation, i.e.,  $i = n$ , and put the dummy variable first in the regression:

$$(24.3.4) \quad \begin{bmatrix} \mathbf{y}(n) \\ \mathbf{y}_n \end{bmatrix} = \begin{bmatrix} \mathbf{o} & \mathbf{X}(n) \\ 1 & \mathbf{x}_n^\top \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \begin{bmatrix} \hat{\varepsilon}(i) \\ \hat{\varepsilon}_n \end{bmatrix} \quad \text{or} \quad \mathbf{y} = \begin{bmatrix} \mathbf{e}_n & \mathbf{X} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \varepsilon$$

• a. 2 points With the definition  $\mathbf{X}_1 = \mathbf{e}_n = \begin{bmatrix} \mathbf{o} \\ 1 \end{bmatrix}$ , write  $\mathbf{M}_1 = \mathbf{I} - \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$  as a  $2 \times 2$  partitioned matrix.

ANSWER.

$$(24.3.5) \quad \mathbf{M}_1 = \begin{bmatrix} \mathbf{I} & \mathbf{o} \\ \mathbf{o}^\top & 1 \end{bmatrix} - \begin{bmatrix} \mathbf{o} \\ 1 \end{bmatrix} \begin{bmatrix} \mathbf{o}^\top & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{o} \\ \mathbf{o}^\top & 0 \end{bmatrix}; \quad \begin{bmatrix} \mathbf{I} & \mathbf{o} \\ \mathbf{o}^\top & 0 \end{bmatrix} \begin{bmatrix} z(i) \\ z_i \end{bmatrix} = \begin{bmatrix} z(i) \\ 0 \end{bmatrix}$$

i.e.,  $\mathbf{M}_1$  simply annuls the last element. □

• b. 2 points Either show mathematically, perhaps by evaluating  $(\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{y}$  or give a good heuristic argument (as [DM93] do), that regressing  $\mathbf{M}_1 \mathbf{y}$  on  $\mathbf{M}_1 \mathbf{X}$  gives the same parameter estimate as regressing  $\mathbf{y}$  on  $\mathbf{X}$  with the  $n$ th observation dropped.

ANSWER. (23.0.2) reads here

$$(24.3.6) \quad \begin{bmatrix} \mathbf{y}(n) \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{X}(n) \\ \mathbf{o}^\top \end{bmatrix} \hat{\beta}(i) + \begin{bmatrix} \hat{\varepsilon}(i) \\ 0 \end{bmatrix}$$

in other words, the estimate of  $\beta$  is indeed  $\hat{\beta}(i)$ , and the first  $n - 1$  elements of the residual indeed the residuals one gets in the regression without the  $i$ th observation. This is so ugly because the singularity shows here in the zeros of the last row, usually it does not show so much. But the way one also sees that it gives zero as the last residual, and this is what one needs to know!

To have a mathematical proof that the last row with zeros does not affect the estimate, evaluate (23.0.3)

$$\begin{aligned} \hat{\beta}_2 &= (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{y} \\ &= \left( \begin{bmatrix} \mathbf{X}(n)^\top & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{o} \\ \mathbf{o}^\top & 0 \end{bmatrix} \begin{bmatrix} \mathbf{X}(n) \\ \mathbf{x}_n^\top \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{X}(n)^\top & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{o} \\ \mathbf{o}^\top & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y}(n) \\ \mathbf{y}_n \end{bmatrix} \\ &= (\mathbf{X}(n)^\top \mathbf{X}(n))^{-1} \mathbf{X}(n)^\top \mathbf{y}(n) = \hat{\beta}(n) \end{aligned}$$

• c. 2 points Use the fact that the residuals in the regression of  $\mathbf{M}_1 \mathbf{y}$  on  $\mathbf{M}_1 \mathbf{X}$  are the same as the residuals in the full regression (24.3.4) to show that  $\hat{\alpha}$  is the  $i$ th predictive residual.

ANSWER.  $\hat{\alpha}$  is obtained from that last row, which reads  $\mathbf{y}_n = \hat{\alpha} + \mathbf{x}_n^\top \hat{\beta}(i)$ , i.e.,  $\hat{\alpha}$  is the predictive residual.

• d. 2 points Use (23.0.3) with  $\mathbf{X}_1$  and  $\mathbf{X}_2$  interchanged to get a formula for  $\hat{\alpha}$ .

ANSWER.  $\hat{\alpha} = (\mathbf{X}_1^\top \mathbf{M} \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{M} \mathbf{y} = \frac{1}{m_{nn}} \hat{\varepsilon}_n = \frac{1}{1 - h_{nn}} \hat{\varepsilon}_n$ , here  $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

• e. 2 points From (23.0.4) follows that also  $\hat{\beta}_2 = (\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top (\mathbf{y} - \mathbf{X}_1 \hat{\beta}_1)$ . Use this to prove

$$(24.3.7) \quad \hat{\beta} - \hat{\beta}(i) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \hat{\varepsilon}_i \frac{1}{1 - h_{ii}}$$

which is [DM93, equation (1.40) on p. 33].

ANSWER. For this we also need to show that one gets the right  $\hat{\beta}(i)$  if one regresses  $\mathbf{y} - \mathbf{e}_n \hat{\varepsilon}_n$  on  $\mathbf{X}$ . In other words,  $\hat{\beta}(n) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{e}_n \hat{\varepsilon}_n)$ , which is exactly (25.4.1).

## CHAPTER 25

## Regression Diagnostics

“Regression Diagnostics” can either concentrate on observations or on variables. Regarding observations, it looks for outliers or influential data in the dataset. Regarding variables, it checks whether there are highly collinear variables, or it keeps track of how much each variable contributes to the  $MSE$  of the regression. Collinearity is discussed in [DM93, 6.3] and [Gre97, 9.2]. Regression diagnostics needs five to ten times more computer resources than the regression itself, and often relies on graphics, therefore it has only recently become part of the standard procedures.

PROBLEM 316. 1 point Define multicollinearity.

- a. 2 points What are the symptoms of multicollinearity?
- b. 2 points How can one detect multicollinearity?
- c. 2 points How can one remedy multicollinearity?

## 25.1. Missing Observations

First case: data on  $\mathbf{y}$  are missing. If you use a least squares predictor then this will not give any change in the estimates and although the computer will think it is more efficient it isn't.

What other schemes are there? Filling in the missing  $\mathbf{y}$  by the arithmetic mean of the observed  $\mathbf{y}$  does not give an unbiased estimator.

General conclusion: in a single-equation context, filling in missing  $\mathbf{y}$  not a good idea.

Now missing values in the  $\mathbf{X}$ -matrix.

If there is only one regressor and a constant term, then the zero order filling in of  $\bar{x}$  “results in no changes and is equivalent with dropping the incomplete data.”

The alternative: filling it with zeros and adding a dummy for the data with missing observation amounts to exactly the same thing.

The only case where filling in missing data makes sense is: if you have multiple regression and you can predict the missing data in the  $\mathbf{X}$  matrix from the other data in the  $\mathbf{X}$  matrix.

## 25.2. Grouped Data

If single observations are replaced by arithmetic means of groups of observations then the error variances vary with the size of the group. If one takes this in consideration, GLS still has good properties, although having the original data is course more efficient.

## 25.3. Influential Observations and Outliers

The following discussion focuses on diagnostics regarding *observations*. To be more precise, we will investigate how each single observation affects the fit established by the other data. (One may also ask how the addition of any *two* observations affects the fit, etc.)

**25.3.1. The “Leverage”.** The  $i$ th diagonal element  $h_{ii}$  of the “hat matrix” is called the “leverage” of the  $i$ th observation. The leverage satisfies the following identity

$$(25.3.1) \quad \hat{\mathbf{y}}_i = (1 - h_{ii})\hat{\mathbf{y}}_i(i) + h_{ii}\mathbf{y}_i$$

$h_{ii}$  is therefore is the weight which  $\mathbf{y}_i$  has in the least squares estimate  $\hat{\mathbf{y}}_i$  of  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$  compared with all other observations, which contribute to  $\hat{\mathbf{y}}_i$  through  $\hat{\mathbf{y}}_i(i)$ . The larger this weight, the more strongly this one observation will influence the estimate of  $\eta_i$  (and if the estimate of  $\eta_i$  is affected, then other parameter estimates may be affected too).

PROBLEM 317. 3 points Explain the meanings of all the terms in equation (25.3.1) and use that equation to explain why  $h_{ii}$  is called the “leverage” of the  $i$ th observation. Is every observation with high leverage also “influential” (in the sense that removal would greatly change the regression estimates)?

ANSWER.  $\hat{\mathbf{y}}_i$  is the fitted value for the  $i$ th observation, i.e., it is the BLUE of  $\eta_i$ , of the expected value of the  $i$ th observation. It is a weighted average of two quantities: the actual observation (which has  $\eta_i$  as expected value), and  $\hat{\mathbf{y}}_i(i)$ , which is the BLUE of  $\eta_i$  based on all the other observations except the  $i$ th. The weight of the  $i$ th observation in this weighted average is called “leverage” of the  $i$ th observation. The sum of all leverages is always  $k$ , the number of parameters in the regression. If the leverage of one individual point is much greater than  $k/n$ , then this point has much more influence on its own fitted value than one should expect just based on the number of observations,

Leverage is not the same as influence; if an observation has high leverage, but by accident the observed value  $\mathbf{y}_i$  is very close to  $\hat{\mathbf{y}}_i(i)$ , then removal of this observation will not change the regression results much. Leverage is potential influence. Leverage does not depend on any of the other observations, one only needs the  $\mathbf{X}$  matrix to compute it.

Those observations whose  $\mathbf{x}$ -values are away from the other observations have high “leverage” and can therefore potentially influence the regression results more than

others.  $h_{ii}$  serves as a measure of this distance. Note that  $h_{ii}$  only depends on the  $\mathbf{X}$ -matrix, not on  $\mathbf{y}$ , i.e., points may have a high leverage but not be influential, because the associated  $\mathbf{y}_i$  blends well into the fit established by the other data. However, regardless of the observed value of  $\mathbf{y}$ , observations with high leverage always affect the covariance matrix of  $\hat{\boldsymbol{\beta}}$ .

$$(25.3.2) \quad h_{ii} = \frac{\det(\mathbf{X}^\top \mathbf{X}) - \det(\mathbf{X}(i)^\top \mathbf{X}(i))}{\det(\mathbf{X}^\top \mathbf{X})},$$

where  $\mathbf{X}(i)$  is the  $\mathbf{X}$ -matrix without the  $i$ th observation.

PROBLEM 318. Prove equation (25.3.2).

ANSWER. Since  $\mathbf{X}^\top(i)\mathbf{X}(i) = \mathbf{X}^\top \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^\top$ , use theorem A.7.3 with  $\mathbf{W} = \mathbf{X}^\top \mathbf{X}$ ,  $\alpha = -1$ , and  $\mathbf{d} = \mathbf{x}_i$ .  $\square$

PROBLEM 319. Prove the following facts about the diagonal elements of the so-called “hat matrix”  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ , which has its name because  $\mathbf{H}\mathbf{y} = \hat{\mathbf{y}}$ , i.e., it puts the hat on  $\mathbf{y}$ .

- a. 1 point  $\mathbf{H}$  is a projection matrix, i.e., it is symmetric and idempotent.

ANSWER. Symmetry follows from the laws for the transposes of products:  $\mathbf{H}^\top = (\mathbf{ABC})^\top = \mathbf{C}^\top \mathbf{B}^\top \mathbf{A}^\top = \mathbf{H}$  where  $\mathbf{A} = \mathbf{X}$ ,  $\mathbf{B} = (\mathbf{X}^\top \mathbf{X})^{-1}$  which is symmetric, and  $\mathbf{C} = \mathbf{X}^\top$ . Idempotency  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .  $\square$

- b. 1 point Prove that a symmetric idempotent matrix is nonnegative definite.

ANSWER. If  $\mathbf{H}$  is symmetric and idempotent, then for arbitrary  $\mathbf{g}$ ,  $\mathbf{g}^\top \mathbf{H} \mathbf{g} = \mathbf{g}^\top \mathbf{H}^\top \mathbf{H} \mathbf{g} = \|\mathbf{H}\mathbf{g}\|^2 \geq 0$ . But  $\mathbf{g}^\top \mathbf{H} \mathbf{g} \geq 0$  for all  $\mathbf{g}$  is the criterion which makes  $\mathbf{H}$  nonnegative definite.  $\square$

- c. 2 points Show that

$$(25.3.3) \quad 0 \leq h_{ii} \leq 1$$

ANSWER. If  $\mathbf{e}_i$  is the vector with a 1 on the  $i$ th place and zeros everywhere else, then  $\mathbf{e}_i^\top \mathbf{H} \mathbf{e}_i = h_{ii}$ . From  $\mathbf{H}$  nonnegative definite follows therefore that  $h_{ii} \geq 0$ .  $h_{ii} \leq 1$  follows because  $\mathbf{I} - \mathbf{H}$  is symmetric and idempotent (and therefore nonnegative definite) as well: it is the projection on the orthogonal complement.  $\square$

- d. 2 points Show: the average value of the  $h_{ii}$  is  $\sum h_{ii}/n = k/n$ , where  $k$  is the number of columns of  $\mathbf{X}$ . (Hint: for this you must compute the trace  $\text{tr } \mathbf{H}$ .)

ANSWER. The average can be written as

$$\frac{1}{n} \text{tr}(\mathbf{H}) = \frac{1}{n} \text{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \frac{1}{n} \text{tr}(\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}) = \frac{1}{n} \text{tr}(\mathbf{I}_k) = \frac{k}{n}.$$

Here we used  $\text{tr } \mathbf{BC} = \text{tr } \mathbf{CB}$  (Theorem A.1.2).  $\square$

- e. 1 point Show that  $\frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top$  is a projection matrix. Here  $\boldsymbol{\iota}$  is the  $n$ -vector of ones.

- f. 2 points Show: If the regression has a constant term, then  $\mathbf{H} - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top$  is a projection matrix.

ANSWER. If  $\boldsymbol{\iota}$ , the vector of ones, is one of the columns of  $\mathbf{X}$  (or a linear combination of these columns), this means there is a vector  $\mathbf{a}$  with  $\boldsymbol{\iota} = \mathbf{X}\mathbf{a}$ . From this follows  $\mathbf{H}\boldsymbol{\iota} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\mathbf{a} = \mathbf{X}\mathbf{a} = \boldsymbol{\iota}$ . One can use this to show that  $\mathbf{H} - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top$  is idempotent:  $(\mathbf{H} - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top)(\mathbf{H} - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top) = \mathbf{H}\mathbf{H} - \mathbf{H}\frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top \mathbf{H} + \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top = \mathbf{H} - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top + \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top = \mathbf{H} - \frac{1}{n} \boldsymbol{\iota} \boldsymbol{\iota}^\top$ .

- g. 1 point Show: If the regression has a constant term, then one can sharpen inequality (25.3.3) to  $1/n \leq h_{ii} \leq 1$ .

ANSWER.  $\mathbf{H} - \boldsymbol{\iota} \boldsymbol{\iota}^\top / n$  is a projection matrix, therefore nonnegative definite, therefore its diagonal elements  $h_{ii} - 1/n$  are nonnegative.

- h. 3 points Why is  $h_{ii}$  called the “leverage” of the  $i$ th observation? To get 3 points, you must give a really good verbal explanation.

ANSWER. Use equation (24.2.12). Effect on any other linear combination of  $\hat{\boldsymbol{\beta}}$  is less than effect on  $\hat{y}_i$ . Distinguish from influence. Leverage depends only on  $\mathbf{X}$  matrix, not on  $\mathbf{y}$ .

$h_{ii}$  is closely related to the test statistic testing whether the  $\mathbf{x}_i$  comes from the same multivariate normal distribution as the other rows of the  $\mathbf{X}$ -matrix. Belskiy, Kuh, and Welsch [BKW80, p. 17] say those observations  $i$  with  $h_{ii} > 2k/n$ , i.e., more than twice the average, should be considered as “leverage points” which might deserve some attention.

### 25.4. Sensitivity of Estimates to Omission of One Observation

The most straightforward approach to sensitivity analysis is to see how the estimates of the parameters of interest are affected if one leaves out the  $i$ th observation. In the case of linear regression, it is not necessary for this to run  $n$  different regressions, but one can derive simple formulas for the changes in the parameters of interest. Interestingly, the various sensitivity measures to be discussed below only depend on the two quantities  $h_{ii}$  and  $\hat{\epsilon}_i$ .

**25.4.1. Changes in the Least Squares Estimate.** Define  $\hat{\boldsymbol{\beta}}(i)$  to be the OLS estimate computed without the  $i$ th observation, and  $\hat{\epsilon}_i(i) = \frac{1}{1-h_{ii}} \hat{\epsilon}_i$  the predictive residual. Then

$$(25.4.1) \quad \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \hat{\epsilon}_i(i)$$

PROBLEM 320. Show (25.4.1) by methods very similar to the proof of (24.2.9).



ANSWER. Here is this brute-force proof, I think from [BKW80]: Let  $\mathbf{y}(i)$  be the  $\mathbf{y}$  vector with the  $i$ th observation deleted. As shown in Problem 308,  $\mathbf{X}^\top(i)\mathbf{y}(i) = \mathbf{X}^\top\mathbf{y} - \mathbf{x}_i\mathbf{y}_i$ . Therefore by (24.2.6)

$$\begin{aligned} \hat{\boldsymbol{\beta}}(i) &= (\mathbf{X}^\top(i)\mathbf{X}(i))^{-1}\mathbf{X}^\top(i)\mathbf{y}(i) = \left( (\mathbf{X}^\top\mathbf{X})^{-1} + \frac{(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}_i^\top(\mathbf{X}^\top\mathbf{X})^{-1}}{1-h_{ii}} \right) (\mathbf{X}^\top\mathbf{y} - \mathbf{x}_i\mathbf{y}_i) \\ &= \hat{\boldsymbol{\beta}} - (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\mathbf{y}_i + \frac{1}{1-h_{ii}}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}_i^\top\hat{\boldsymbol{\beta}} - \frac{h_{ii}}{1-h_{ii}}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\mathbf{y}_i \\ &= \hat{\boldsymbol{\beta}} - \frac{1}{1-h_{ii}}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\mathbf{y}_i + \frac{1}{1-h_{ii}}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}_i^\top\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - \frac{1}{1-h_{ii}}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\hat{\boldsymbol{\epsilon}}_i \end{aligned}$$

□

To understand (25.4.1), note the following fact which is interesting in its own right:  $\hat{\boldsymbol{\beta}}(i)$ , which is defined as the OLS estimator if one drops the  $i$ th observation, can also be obtained as the OLS estimator if one replaces the  $i$ th observation by the prediction of the  $i$ th observation on the basis of all other observations, i.e., by  $\hat{\mathbf{y}}_i(i)$ . Writing  $\mathbf{y}((i))$  for the vector  $\mathbf{y}$  whose  $i$ th observation has been replaced in this way, one obtains

$$(25.4.2) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}; \quad \hat{\boldsymbol{\beta}}(i) = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}((i)).$$

Since  $\mathbf{y} - \mathbf{y}((i)) = \mathbf{e}_i\hat{\boldsymbol{\epsilon}}_i(i)$  and  $\mathbf{x}_i = \mathbf{X}^\top\mathbf{e}_i$  (25.4.1) follows.

The quantities  $h_{ii}$ ,  $\hat{\boldsymbol{\beta}}(i) - \hat{\boldsymbol{\beta}}$ , and  $s^2(i)$  are computed by the R-function `lm.influence`. Compare [CH93, pp. 129–131].

**25.4.2. Scaled Measures of Sensitivity.** In order to assess the sensitivity of the estimate of any linear combination of the elements of  $\boldsymbol{\beta}$ ,  $\phi = \mathbf{t}^\top\boldsymbol{\beta}$ , it makes sense to divide the change in  $\mathbf{t}^\top\hat{\boldsymbol{\beta}}$  due to omission of the  $i$ th observation by the standard deviation of  $\mathbf{t}^\top\hat{\boldsymbol{\beta}}$ , i.e., to look at

$$(25.4.3) \quad \frac{\mathbf{t}^\top(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i))}{\sigma\sqrt{\mathbf{t}^\top(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{t}}}$$

Such a standardization makes it possible to compare the sensitivity of different linear combinations, and to ask: Which linear combination of the elements of  $\hat{\boldsymbol{\beta}}$  is affected most if one drops the  $i$ th observation? Interestingly and, in hindsight, perhaps not surprisingly, the linear combination which is most sensitive to the addition of the  $i$ th observation, is  $\mathbf{t} = \mathbf{x}_i$ .

For a mathematical proof we need the following inequality, which is nothing but the Cauchy-Schwartz inequality in disguise:

THEOREM 25.4.1. *If  $\boldsymbol{\Omega}$  is positive definite symmetric, then*

$$(25.4.4) \quad \max_g \frac{(\mathbf{g}^\top\mathbf{x})^2}{\mathbf{g}^\top\boldsymbol{\Omega}\mathbf{g}} = \mathbf{x}^\top\boldsymbol{\Omega}^{-1}\mathbf{x}.$$

*If the denominator in the fraction on the lefthand side is zero, then  $\mathbf{g} = \mathbf{o}$  and therefore the numerator is necessarily zero as well. In this case, the fraction itself should be considered zero.*

Proof: As in the derivation of the BLUE with nonspherical covariance matrix, pick a nonsingular  $\mathbf{Q}$  with  $\boldsymbol{\Omega} = \mathbf{Q}\mathbf{Q}^\top$ , and define  $\mathbf{P} = \mathbf{Q}^{-1}$ . Then it follows  $\mathbf{P}\boldsymbol{\Omega}\mathbf{P}^\top = \mathbf{I}$ . Define  $\mathbf{y} = \mathbf{P}\mathbf{x}$  and  $\mathbf{h} = \mathbf{Q}^\top\mathbf{g}$ . Then  $\mathbf{h}^\top\mathbf{y} = \mathbf{g}^\top\mathbf{x}$ ,  $\mathbf{h}^\top\mathbf{h} = \mathbf{g}^\top\boldsymbol{\Omega}\mathbf{g}$ , and  $\mathbf{y}^\top\mathbf{y} = \mathbf{x}^\top\boldsymbol{\Omega}^{-1}\mathbf{x}$ . Therefore (25.4.4) follows from the Cauchy-Schwartz inequality  $(\mathbf{h}^\top\mathbf{y})^2 \leq (\mathbf{h}^\top\mathbf{h})(\mathbf{y}^\top\mathbf{y})$ .

Using Theorem 25.4.1 and equation (25.4.1) one obtains

$$(25.4.5) \quad \max_{\mathbf{t}} \frac{(\mathbf{t}^\top(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i)))^2}{\sigma^2\mathbf{t}^\top(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{t}} = \frac{1}{\sigma^2}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i))^\top\mathbf{X}^\top\mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i)) = \frac{1}{\sigma^2}\mathbf{x}_i^\top(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{x}_i\hat{\boldsymbol{\epsilon}}_i^2(i) = \frac{h_{ii}}{\sigma^2}\hat{\boldsymbol{\epsilon}}_i^2(i)$$

Now we will show that the linear combination which attains this maximum, which is most sensitive to the addition of the  $i$ th observation, is  $\mathbf{t} = \mathbf{x}_i$ . If one premultiplies (25.4.1) by  $\mathbf{x}_i^\top$  one obtains

$$(25.4.6) \quad \hat{y}_i - \hat{y}_i(i) = \mathbf{x}_i^\top\hat{\boldsymbol{\beta}} - \mathbf{x}_i^\top\hat{\boldsymbol{\beta}}(i) = \frac{h_{ii}}{1-h_{ii}}\hat{\boldsymbol{\epsilon}}_i = h_{ii}\hat{\boldsymbol{\epsilon}}_i(i)$$

If one divides (25.4.6) by the standard deviation of  $\hat{y}_i$ , i.e., if one applies the construction (25.4.3), one obtains

$$(25.4.7) \quad \frac{\hat{y}_i - \hat{y}_i(i)}{\sigma\sqrt{h_{ii}}} = \frac{\sqrt{h_{ii}}}{\sigma}\hat{\boldsymbol{\epsilon}}_i(i) = \frac{\sqrt{h_{ii}}}{\sigma(1-h_{ii})}\hat{\boldsymbol{\epsilon}}_i$$

If  $\hat{y}_i$  changes only little (compared with the standard deviation of  $\hat{y}_i$ ) if the observation is removed, then no other linear combination of the elements of  $\hat{\boldsymbol{\beta}}$  will be affected much by the omission of this observation either.

The righthand side of (25.4.7), with  $\sigma$  estimated by  $s(i)$ , is called by [BKW80] and many others DFFITS (which stands for DiFference in FIT, Standardized). If one takes its square, divides it by  $k$ , and estimates  $\sigma^2$  by  $s^2$  (which is more consistent than using  $s^2(i)$ , since one standardizes by the standard deviation of  $\mathbf{t}^\top\hat{\boldsymbol{\beta}}$  and not by that of  $\mathbf{t}^\top\hat{\boldsymbol{\beta}}(i)$ ), one obtains Cook's distance [Coo77]. (25.4.5) gives an equation for Cook's distance in terms of  $\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i)$ :

$$(25.4.8) \quad \text{Cook's distance} = \frac{(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i))^\top\mathbf{X}^\top\mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}(i))}{ks^2} = \frac{h_{ii}}{ks^2}\hat{\boldsymbol{\epsilon}}_i^2(i) = \frac{h_{ii}}{ks^2(1-h_{ii})^2}\hat{\boldsymbol{\epsilon}}_i^2(i)$$

PROBLEM 321. *Can you think of a situation in which an observation has a small residual but a large "influence" as measured by Cook's distance?*

ANSWER. Assume “all observations are clustered near each other while the solitary odd observation lies a way out” as Kmenta wrote in [Kme86, p. 426]. If the observation happens to lie on the regression line, then it can be discovered by its influence on the variance-covariance matrix (25.3.2), i.e., in this case only the  $h_{ii}$  count.  $\square$

PROBLEM 322. *The following is the example given in [Coo77]. In R, the command `data(longley)` makes the data frame `longley` available, which has the famous Longley-data, a standard example for a highly multicollinear dataset. These data are also available on the web at [www.econ.utah.edu/ehrbbar/data/longley.txt](http://www.econ.utah.edu/ehrbbar/data/longley.txt). `attach(longley)` makes the individual variables available as R-objects.*

• a. *3 points Look at the data in a scatterplot matrix and explain what you see. Later we will see that one of the observations is in the regression much more influential than the rest. Can you see from the scatterplot matrix which observation that might be?*

ANSWER. In linux, you first have to give the command `x11()` in order to make the graphics window available. In windows, this is not necessary. It is important to display the data in a reasonable order, therefore instead of `pairs(longley)` you should do something like `attach(longley)` and then `pairs(cbind(Year, Population, Employed, Unemployed, Armed.Forces, GNP, GNP.deflator))`. Put `Year` first, so that all variables are plotted against `Year` on the horizontal axis.

Population vs. year is a very smooth line.

Population vs GNP also quite smooth.

You see the huge increase in the armed forced in 1951 due to the Korean War, which led to a (temporary) drop in unemployment and a (not so temporary) jump in the GNP deflator.

Otherwise the unemployed show the stop-and-go scenario of the fifties.

unemployed is not correlated with anything.

One should expect a strong negative correlation between employed and unemployed, but this is not the case.  $\square$

• b. *4 points Run a regression of the model `Employed ~ GNP.deflator + GNP + Unemployed + Armed.Forces + Population + Year` and discuss the result.*

ANSWER. To fit a regression run `longley.fit <- lm(Employed ~ GNP + Unemployed + Armed.Forces + Population + Year)`. You can see the regression results by typing `summary(longley.fit)`.

Armed forces and unemployed are significant and have negative sign, as expected.

GNP and Population are insignificant and have negative sign too, this is not expected. GNP, Population and Year are highly collinear.  $\square$

• c. *3 points Make plots of the ordinary residuals and the standardized residuals against time. How do they differ? In R, the commands are `plot(Year, residuals(longley.fit), type="h", ylab="Ordinary Residuals in Longley Regression")`. In order to get the next plot in a different graphics window, so that you can compare them, do now either `x11()` in linux or `windows()` in windows, and then `plot(Year, rstandard(longley.fit), type="h", ylab="Standardized Residuals in Longley Regression")`.*

ANSWER. You see that the standardized residuals at the edge of the dataset are bigger than the ordinary residuals. The datapoints at the edge are better able to attract the regression plane than those in the middle, therefore the ordinary residuals are “too small.” Standardization corrects for this.

• d. *4 points Make plots of the predictive residuals. Apparently there is no speed of convergence command in R to do this, therefore you should use formula (24.2.9). Also plot standardized predictive residuals, and compare them.*

ANSWER. The predictive residuals are `plot(Year, residuals(longley.fit)/(1-hatvalues(longley.fit)), type="h", ylab="Predictive Residuals in Longley Regression")`. The standardized predictive residuals are often called studentized residuals, `plot(Year, rstudent(longley.fit), type="h", ylab="Standardized predictive Residuals in Longley Regression")`.

A comparison shows an opposite effect as with the ordinary residuals: the predictive residuals at the edge of the dataset are too large, and standardization corrects this.

Specific results: standardized predictive residual in 1950 smaller than that in 1962, but predictive residual in 1950 is very close to 1962.

standardized predictive residual in 1951 smaller than that in 1956, but predictive residual in 1951 is larger than in 1956.

Largest predictive residual is 1951, but largest standardized predictive residual is 1956.

• e. *3 points Make a plot of the leverage, i.e., the  $h_{ii}$ -values, using `plot(Year, hatvalues(longley.fit), type="h", ylab="Leverage in Longley Regression")` and explain what leverage means.*

• f. *3 points One observation is much more influential than the others; which is it? First look at the plots for the residuals, then look also at the plot for leverage and try to guess which is the most influential observation. Then do it the right way. Can you give reasons based on your prior knowledge about the time period involved why an observation in that year might be influential?*

ANSWER. The “right” way is to use Cook’s distance: `plot(Year, cooks.distance(longley.fit), type="h", ylab="Cook’s Distance in Longley Regression")`

One sees that 1951 towers above all others. It does not have highest leverage, but it is second-highest, and a bigger residual than the point with the highest leverage.

1951 has the largest distance of .61. The second largest is the last observation in the dataset, 1962, with a distance of .47, and the others have .24 or less. Cook says: removal of 1951 point would move the least squares estimate to the edge of a 35% confidence region around  $\hat{\beta}$ . This point is probably so influential because 1951 was the first full year of the Korean war. One would not be able to detect this point from the ordinary residuals, standardized or not! The predictive residuals are a little better; their maximum is at 1951, but several other residuals are almost as large. 1951 is so influential because it has an extremely high hat-value, and one of the highest values for ordinary residuals!

*At the end don’t forget to `detach(longley)` if you have attached it before.*

**25.4.3. Changes in the Sum of Squared Errors.** For the computation of  $s^2(i)$  from the regression results one can take advantage of the following simple relationship between the SSE for the regression with and without the  $i$ th observation:

$$(25.4.9) \quad \text{SSE} - \text{SSE}(i) = \frac{\hat{\epsilon}_i^2}{1 - h_{ii}}$$

PROBLEM 323. Use (25.4.9) to derive the following formula for  $s^2(i)$ :

$$(25.4.10) \quad s^2(i) = \frac{1}{n - k - 1} \left( (n - k)s^2 - \frac{\hat{\epsilon}_i^2}{1 - h_{ii}} \right)$$

ANSWER. This merely involves re-writing SSE and SSE( $i$ ) in terms of  $s^2$  and  $s^2(i)$ .

$$(25.4.11) \quad s^2(i) = \frac{\text{SSE}(i)}{n - 1 - k} = \frac{1}{n - k - 1} \left( \text{SSE} - \frac{\hat{\epsilon}_i^2}{1 - h_{ii}} \right)$$

□

*Proof of equation (25.4.9):*

$$\begin{aligned} \text{SSE}(i) &= \sum_{j: j \neq i} (\mathbf{y}_j - \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}(i))^2 = \sum_{j: j \neq i} (\mathbf{y}_j - \mathbf{x}_j^\top \hat{\boldsymbol{\beta}} - \mathbf{x}_j^\top (\hat{\boldsymbol{\beta}}(i) - \hat{\boldsymbol{\beta}}))^2 \\ &= \sum_{j: j \neq i} \left( \hat{\epsilon}_j + \frac{h_{ji}}{1 - h_{ii}} \hat{\epsilon}_i \right)^2 \\ &= \sum_j \left( \hat{\epsilon}_j + \frac{h_{ji}}{1 - h_{ii}} \hat{\epsilon}_i \right)^2 - \left( \frac{1}{1 - h_{ii}} \hat{\epsilon}_i \right)^2 \\ &= \sum_j \hat{\epsilon}_j^2 + \frac{2\hat{\epsilon}_i}{1 - h_{ii}} \sum_j h_{ij} \hat{\epsilon}_j + \left( \frac{\hat{\epsilon}_i}{1 - h_{ii}} \right)^2 \sum_j h_{ji}^2 - \left( \frac{\hat{\epsilon}_i}{1 - h_{ii}} \right)^2 \end{aligned}$$

In the last line the first term is SSE. The second term is zero because  $\mathbf{H}\hat{\boldsymbol{\epsilon}} = \mathbf{o}$ . Furthermore,  $h_{ii} = \sum_j h_{ji}^2$  because  $\mathbf{H}$  is symmetric and idempotent, therefore the sum of the last two items is  $-\hat{\epsilon}_i^2/(1 - h_{ii})$ .

Note that every single relationship we have derived so far is a function of  $\hat{\epsilon}_i$  and  $h_{ii}$ .

PROBLEM 324. 3 points What are the main concepts used in modern “Regression Diagnostics”? Can it be characterized to be a careful look at the residuals, or does it have elements which cannot be inferred from the residuals alone?

ANSWER. Leverage (sometimes it is called “potential”) is something which cannot be inferred from the residuals, it does not depend on  $\mathbf{y}$  at all. □

PROBLEM 325. An observation in a linear regression model is “influential” if its omission causes large changes to the regression results. Discuss how you would ascertain in practice whether a given observation is influential or not.

• a. What is meant by leverage? Does high leverage necessarily imply that an observation is influential?

ANSWER. Leverage is potential influence. It only depends of  $\mathbf{X}$ , not on  $\mathbf{y}$ . It is the distance of the observation from the center of gravity of all observations. Whether this is actual influence depends on the  $\mathbf{y}$ -values.

• b. How are the concepts of leverage and influence affected by sample size?

• c. What steps would you take when alerted to the presence of an influential observation?

ANSWER. Make sure you know whether the results you rely on are affected if that influential observation is dropped. Try to find out why this observation is influential (e.g. in the Longley data the observations in the year when the Korean War started are influential).

• d. What is a “predictive residual” and how does it differ from an ordinary residual?

• e. Discuss situations in which one would want to deal with the “predictive residuals” rather than the ordinary residuals, and situations in which one would want standardized residuals versus situations in which it would be preferable to have unstandardized residuals.

PROBLEM 326. 6 points Describe what you would do to ascertain that a regression you ran is correctly specified?

ANSWER. Economic theory behind that regression, size and sign of coefficients, plot residuals versus predicted values, time, and every independent variable, run all tests: F-test, t-tests, DW, portmanteau test, forecasting, multicollinearity, influence statistics, overfitting to see if other variables are significant, try to defeat the result by using alternative variables, divide time period into subperiods in order to see if parameters are constant over time, pre-test specification assumptions.

## CHAPTER 26

## Asymptotic Properties of the OLS Estimator

A much more detailed treatment of the contents of this chapter can be found in [DM93, Chapters 4 and 5].

Here we are concerned with the consistency of the OLS estimator for large samples. In other words, we assume that our regression model can be extended to encompass an arbitrary number of observations. First we assume that the regressors are nonstochastic, and we will make the following assumption:

$$(26.0.12) \quad \mathbf{Q} = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^\top \mathbf{X} \text{ exists and is nonsingular.}$$

Two examples where this is not the case. Look at the model  $\mathbf{y}_t = \alpha + \beta t + \varepsilon_t$ . Here

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ \vdots & \vdots \\ 1 & n \end{bmatrix}. \text{ Therefore } \mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 1+1+1+\cdots+1 & 1+2+3+\cdots+n \\ 1+2+3+\cdots+n & 1+4+9+\cdots+n^2 \end{bmatrix} =$$

$\begin{bmatrix} n & n(n+1)/2 \\ n(n+1)/2 & n(n+1)(2n+1)/6 \end{bmatrix}$ , and  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} \rightarrow \begin{bmatrix} 1 & \infty \\ \infty & \infty \end{bmatrix}$ . Here the assumption (26.0.12) does not hold, but one can still prove consistency and asymptotic normality, the estimators converge even faster than in the usual case.

The other example is the model  $\mathbf{y}_t = \alpha + \beta \lambda^t + \varepsilon_t$  with a known  $\lambda$  with  $-1 < \lambda < 1$ . Here

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} &= \begin{bmatrix} 1+1+\cdots+1 & \lambda+\lambda^2+\cdots+\lambda^n \\ \lambda+\lambda^2+\cdots+\lambda^n & \lambda^2+\lambda^4+\cdots+\lambda^{2n} \end{bmatrix} = \\ &= \begin{bmatrix} n & (\lambda-\lambda^{n+1})/(1-\lambda) \\ (\lambda-\lambda^{n+1})/(1-\lambda) & (\lambda^2-\lambda^{2n+2})/(1-\lambda^2) \end{bmatrix}. \end{aligned}$$

Therefore  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} \rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ , which is singular. In this case, a consistent estimate of  $\lambda$  does not exist: future observations depend on  $\lambda$  so little that even with infinitely many observations there is not enough information to get the precise value of  $\lambda$ .

We will show that under assumption (26.0.12),  $\hat{\boldsymbol{\beta}}$  and  $s^2$  are *consistent*. However, this assumption is really too strong for consistency. A weaker set of assumptions is the Grenander conditions, see [Gre97, p. 275]. To write down the Grenander conditions, remember that presently  $\mathbf{X}$  depends on  $n$  (in that we only look at the first  $n$  elements of  $\mathbf{y}$  and first  $n$  rows of  $\mathbf{X}$ ), therefore also the column vectors  $\mathbf{x}_j$  also depend on  $n$  (although we are not indicating this here). Therefore  $\mathbf{x}_j^\top \mathbf{x}_j$  depends on  $n$  as well, and we will make this dependency explicit by writing  $\mathbf{x}_j^\top \mathbf{x}_j = d_{nj}^2$ . The first Grenander condition is  $\lim_{n \rightarrow \infty} d_{nj}^2 = +\infty$  for all  $j$ . Second: for all  $i$  and  $k$ ,  $\lim_{n \rightarrow \infty} \max_{i=1 \dots n} x_{ij}/d_{nj}^2 = 0$  (here is a typo in Greene, he leaves the  $\max$  out). Third: Sample correlation matrix of the columns of  $\mathbf{X}$  minus the constant term converges to a nonsingular matrix.

Consistency means that the *probability limit* of the estimates converges towards the true value. For  $\hat{\boldsymbol{\beta}}$  this can be written as  $\text{plim}_{n \rightarrow \infty} \hat{\boldsymbol{\beta}}_n = \boldsymbol{\beta}$ . This means that for all  $\varepsilon > 0$  follows  $\lim_{n \rightarrow \infty} \Pr[|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}| \leq \varepsilon] = 1$ .

The probability limit is one of several concepts of limits used in probability theory. We will need the following properties of the plim here:

(1) For nonrandom magnitudes, the probability limit is equal to the ordinary limit.

(2) It satisfies the Slutsky theorem, that for a continuous function  $g$ ,

$$(26.0.13) \quad \text{plim } g(\mathbf{z}) = g(\text{plim}(\mathbf{z})).$$

(3) If the  $\mathcal{MSE}$ -matrix of an estimator converges towards the null matrix, then the estimator is consistent.

(4) Kinchine's theorem: the sample mean of an i.i.d. distribution is a consistent estimate of the population mean, even if the distribution does not have a population variance.

## 26.1. Consistency of the OLS estimator

For the proof of consistency of the OLS estimators  $\hat{\boldsymbol{\beta}}$  and of  $s^2$  we need the following result:

$$(26.1.1) \quad \text{plim } \frac{1}{n} \mathbf{X}^\top \boldsymbol{\varepsilon} = \mathbf{0}.$$

I.e., the true  $\boldsymbol{\varepsilon}$  is asymptotically orthogonal to all columns of  $\mathbf{X}$ . This follows immediately from  $\mathcal{MSE}[\boldsymbol{\varepsilon}; \mathbf{X}^\top \boldsymbol{\varepsilon}/n] = \mathcal{E}[\mathbf{X}^\top \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \mathbf{X}/n^2] = \sigma^2 \mathbf{X}^\top \mathbf{X}/n^2$ , which converges towards  $\mathbf{O}$ .

In order to prove consistency of  $\hat{\boldsymbol{\beta}}$  and  $s^2$ , transform the formulas for  $\hat{\boldsymbol{\beta}}$  and  $s^2$  in such a way that they are written as continuous functions of terms each of which

converges for  $n \rightarrow \infty$ , and then apply Slutsky's theorem. Write  $\hat{\boldsymbol{\beta}}$  as

$$(26.1.2) \quad \hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} = \boldsymbol{\beta} + \left( \frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \frac{\mathbf{X}^\top \boldsymbol{\varepsilon}}{n}$$

$$(26.1.3) \quad \text{plim } \hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + \lim \left( \frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \text{plim } \frac{\mathbf{X}^\top \boldsymbol{\varepsilon}}{n}$$

$$(26.1.4) \quad = \boldsymbol{\beta} + \mathbf{Q}^{-1} \mathbf{o} = \boldsymbol{\beta}.$$

Let's look at the geometry of this when there is only one explanatory variable. The specification is therefore  $\mathbf{y} = \mathbf{x}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ . The assumption is that  $\boldsymbol{\varepsilon}$  is asymptotically orthogonal to  $\mathbf{x}$ . In small samples, it only happens by sheer accident with probability 0 that  $\boldsymbol{\varepsilon}$  is orthogonal to  $\mathbf{x}$ . Only  $\hat{\boldsymbol{\varepsilon}}$  is. But now let's assume the sample grows larger, i.e., the vectors  $\mathbf{y}$  and  $\mathbf{x}$  become very high-dimensional observation vectors, i.e. we are drawing here a two-dimensional subspace out of a very high-dimensional space. As more and more data are added, the observation vectors also become longer and longer. But if we divide each vector by  $\sqrt{n}$ , then the lengths of these normalized lengths stabilize. The squared length of the vector  $\boldsymbol{\varepsilon}/\sqrt{n}$  has the plim of  $\sigma^2$ . Furthermore, assumption (26.0.12) means in our case that  $\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{x}^\top \mathbf{x}$  exists and is nonsingular. This is the squared length of  $\frac{1}{\sqrt{n}} \mathbf{x}$ . I.e., if we normalize the vectors by dividing them by  $\sqrt{n}$ , then they do not get longer but converge towards a finite length. And the result (26.1.1)  $\text{plim } \frac{1}{n} \mathbf{x}^\top \boldsymbol{\varepsilon} = 0$  means now that with this normalization,  $\boldsymbol{\varepsilon}/\sqrt{n}$  becomes more and more orthogonal to  $\mathbf{x}/\sqrt{n}$ . I.e., if  $n$  is large enough, asymptotically, not only  $\hat{\boldsymbol{\varepsilon}}$  but also the true  $\boldsymbol{\varepsilon}$  is orthogonal to  $\mathbf{x}$ , and this means that asymptotically  $\hat{\boldsymbol{\beta}}$  converges towards the true  $\boldsymbol{\beta}$ .

For the proof of consistency of  $s^2$  we need, among others, that  $\text{plim } \frac{\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}}{n} = \sigma^2$ , which is a consequence of Kinchine's theorem. Since  $\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon}^\top \mathbf{M} \boldsymbol{\varepsilon}$  it follows

$$\begin{aligned} \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n-k} &= \frac{n}{n-k} \boldsymbol{\varepsilon}^\top \left( \frac{\mathbf{I}}{n} - \frac{\mathbf{X}}{n} \left( \frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \frac{\mathbf{X}^\top}{n} \right) \boldsymbol{\varepsilon} = \\ &= \frac{n}{n-k} \left( \frac{\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}}{n} - \frac{\boldsymbol{\varepsilon}^\top \mathbf{X}}{n} \left( \frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \frac{\mathbf{X}^\top \boldsymbol{\varepsilon}}{n} \right) \rightarrow 1 \cdot \left( \sigma^2 - \mathbf{o}^\top \mathbf{Q}^{-1} \mathbf{o} \right). \end{aligned}$$

## 26.2. Asymptotic Normality of the Least Squares Estimator

To show asymptotic normality of an estimator, multiply the sampling error by  $\sqrt{n}$ , so that the variance is stabilized.

We have seen  $\text{plim } \frac{1}{n} \mathbf{X}^\top \boldsymbol{\varepsilon} = \mathbf{o}$ . Now look at  $\frac{1}{\sqrt{n}} \mathbf{X}^\top \boldsymbol{\varepsilon}_n$ . Its mean is  $\mathbf{o}$  and its covariance matrix  $\sigma^2 \frac{\mathbf{X}^\top \mathbf{X}}{n}$ . Shape of distribution, due to a variant of the Central Limit Theorem, is asymptotically normal:  $\frac{1}{\sqrt{n}} \mathbf{X}^\top \boldsymbol{\varepsilon}_n \rightarrow \text{N}(\mathbf{o}, \sigma^2 \mathbf{Q})$ . (Here the convergence is convergence in distribution.)

We can write  $\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) = \left( \frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \left( \frac{1}{\sqrt{n}} \mathbf{X}^\top \boldsymbol{\varepsilon}_n \right)$ . Therefore its limiting covariance matrix is  $\mathbf{Q}^{-1} \sigma^2 \mathbf{Q} \mathbf{Q}^{-1} = \sigma^2 \mathbf{Q}^{-1}$ . Therefore  $\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \rightarrow \text{N}(\mathbf{o}, \sigma^2 \mathbf{Q}^{-1})$  in distribution. One can also say: the asymptotic distribution of  $\hat{\boldsymbol{\beta}}$  is  $\text{N}(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$ .

From this follows  $\sqrt{n}(\mathbf{R}\hat{\boldsymbol{\beta}}_n - \mathbf{R}\boldsymbol{\beta}) \rightarrow \text{N}(\mathbf{o}, \sigma^2 \mathbf{R}\mathbf{Q}^{-1}\mathbf{R}^\top)$ , and therefore

$$(26.2.1) \quad n(\mathbf{R}\hat{\boldsymbol{\beta}}_n - \mathbf{R}\boldsymbol{\beta})(\mathbf{R}\mathbf{Q}^{-1}\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}}_n - \mathbf{R}\boldsymbol{\beta}) \rightarrow \sigma^2 \chi_i^2.$$

Divide by  $s^2$  and replace in the limiting case  $\mathbf{Q}$  by  $\mathbf{X}^\top \mathbf{X}/n$  and  $s^2$  by  $\sigma^2$  to get

$$(26.2.2) \quad \frac{(\mathbf{R}\hat{\boldsymbol{\beta}}_n - \mathbf{R}\boldsymbol{\beta})(\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}}_n - \mathbf{R}\boldsymbol{\beta})}{s^2} \rightarrow \chi_i^2$$

in distribution. All this is not a proof; the point is that in the denominator, the distribution is divided by the increasingly bigger number  $n - k$ , while in the numerator, it is divided by the constant  $i$ ; therefore asymptotically the denominator can be considered 1.

The central limit theorems only say that for  $n \rightarrow \infty$  these converge towards  $\chi^2$ , which is asymptotically equal to the F distribution. It is easily possible that before one gets to the limit, the F-distribution is better.

**PROBLEM 327.** *Are the residuals  $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  asymptotically normally distributed?*

**ANSWER.** Only if the disturbances are normal, otherwise of course not! We can show that  $\sqrt{n}(\boldsymbol{\varepsilon} - \hat{\boldsymbol{\varepsilon}}) = \sqrt{n}\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim \text{N}(\mathbf{o}, \sigma^2 \mathbf{X}\mathbf{Q}\mathbf{X}^\top)$ .

Now these results also go through if one has stochastic regressors. [Gre97, 6.7] shows that the above condition (26.0.12) with the lim replaced by plim holds if  $\boldsymbol{\varepsilon}_i$  and  $\mathbf{x}_i$  are an i.i.d. sequence of random variables.

**PROBLEM 328.** *2 points In the regression model with random regressors  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , you only know that  $\text{plim } \frac{1}{n} \mathbf{X}^\top \mathbf{X} = \mathbf{Q}$  is a nonsingular matrix, and  $\text{plim } \frac{1}{n} \mathbf{X}^\top \boldsymbol{\varepsilon} = \mathbf{o}$ . Using these two conditions, show that the OLS estimate is consistent.*

**ANSWER.**  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}$  due to (18.0.7), and

$$\text{plim}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} = \text{plim} \left( \frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \frac{\mathbf{X}^\top \boldsymbol{\varepsilon}}{n} = \mathbf{Q} \mathbf{o} = \mathbf{o}.$$

## CHAPTER 27

## Least Squares as the Normal Maximum Likelihood Estimate

Now assume  $\boldsymbol{\varepsilon}$  is multivariate normal. We will show that in this case the OLS estimator  $\hat{\boldsymbol{\beta}}$  is at the same time the Maximum Likelihood Estimator. For this we need to write down the density function of  $\mathbf{y}$ . First look at one  $\mathbf{y}_t$  which is  $\mathbf{y}_t \sim$

$\mathbf{N}(\mathbf{x}_t^\top \boldsymbol{\beta}, \sigma^2)$ , where  $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix}$ , i.e.,  $\mathbf{x}_t$  is the  $t$ th row of  $\mathbf{X}$ . It is written as a column vector, since we follow the “column vector convention.” The (marginal) density function for this one observation is

$$(27.0.3) \quad f_{\mathbf{y}_t}(y_t) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_t - \mathbf{x}_t^\top \boldsymbol{\beta})^2 / 2\sigma^2}.$$

Since the  $\mathbf{y}_i$  are stochastically independent, their joint density function is the product, which can be written as

$$(27.0.4) \quad f_{\mathbf{y}}(\mathbf{y}) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right).$$

To compute the maximum likelihood estimator, it is advantageous to start with the log likelihood function:

$$(27.0.5) \quad \log f_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

Assume for a moment that  $\sigma^2$  is known. Then the MLE of  $\boldsymbol{\beta}$  is clearly equal to the OLS  $\hat{\boldsymbol{\beta}}$ . Since  $\hat{\boldsymbol{\beta}}$  does not depend on  $\sigma^2$ , it is also the maximum likelihood estimate when  $\sigma^2$  is unknown.  $\hat{\boldsymbol{\beta}}$  is a linear function of  $\mathbf{y}$ . Linear transformations of normal variables are normal. Normal distributions are characterized by their mean vector and covariance matrix. The distribution of the MLE of  $\boldsymbol{\beta}$  is therefore  $\hat{\boldsymbol{\beta}} \sim \mathbf{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$ .

If we replace  $\boldsymbol{\beta}$  in the log likelihood function (27.0.5) by  $\hat{\boldsymbol{\beta}}$ , we get what is called the log likelihood function with  $\boldsymbol{\beta}$  “concentrated out.”

$$(27.0.6) \quad \log f_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\beta} = \hat{\boldsymbol{\beta}}, \sigma^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

One gets the maximum likelihood estimate of  $\sigma^2$  by maximizing this “concentrated log likelihood function.” Taking the derivative with respect to  $\sigma^2$  (consider  $\sigma^2$  the name of a variable, not the square of another variable), one gets

$$(27.0.7) \quad \frac{\partial}{\partial \sigma^2} \log f_{\mathbf{y}}(\mathbf{y}; \hat{\boldsymbol{\beta}}) = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

Setting this zero gives

$$(27.0.8) \quad \hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n} = \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n}.$$

This is a scalar multiple of the unbiased estimate  $s^2 = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} / (n - k)$  which we had earlier.

Let’s look at the distribution of  $s^2$  (from which that of its scalar multiples follows easily). It is a quadratic form in a normal variable. Such quadratic forms very often have  $\chi^2$  distributions.

Now recall equation 7.4.9 characterizing all the quadratic forms of multivariate normal variables that are  $\chi^2$ ’s. Here it is again: Assume  $\mathbf{y}$  is a multivariate normal vector random variable with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\sigma^2 \boldsymbol{\Psi}$ , and  $\boldsymbol{\Omega}$  is a symmetric nonnegative definite matrix. Then  $(\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Omega}(\mathbf{y} - \boldsymbol{\mu}) \sim \sigma^2 \chi_k^2$  iff

$$(27.0.9) \quad \boldsymbol{\Psi} \boldsymbol{\Omega} \boldsymbol{\Psi} \boldsymbol{\Omega} \boldsymbol{\Psi} = \boldsymbol{\Psi} \boldsymbol{\Omega} \boldsymbol{\Psi},$$

and  $k$  is the rank of  $\boldsymbol{\Psi} \boldsymbol{\Omega}$ .

This condition is satisfied in particular if  $\boldsymbol{\Psi} = \mathbf{I}$  (the identity matrix) and  $\boldsymbol{\Omega}^2 = \boldsymbol{\Omega}$ , and this is exactly our situation.

$$(27.0.10) \quad \hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n - k} = \frac{\boldsymbol{\varepsilon}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \boldsymbol{\varepsilon}}{n - k} = \frac{\boldsymbol{\varepsilon}^\top \mathbf{M} \boldsymbol{\varepsilon}}{n - k}$$

where  $\mathbf{M}^2 = \mathbf{M}$  and  $\text{rank } \mathbf{M} = n - k$ . (This last identity because for idempotent matrices,  $\text{rank} = \text{tr}$ , and we computed its  $\text{tr}$  above.) Therefore  $s^2 \sim \sigma^2 \chi_{n-k}^2 / (n - k)$  from which one obtains again unbiasedness, but also that  $\text{var}[s^2] = 2\sigma^4 / (n - k)$  result that one cannot get from mean and variance alone.

**PROBLEM 329.** 4 points Show that, if  $\mathbf{y}$  is normally distributed,  $s^2$  and  $\hat{\boldsymbol{\beta}}$  are independent.

**ANSWER.** We showed in question 246 that  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\varepsilon}}$  are uncorrelated, therefore in the normal case independent, therefore  $\hat{\boldsymbol{\beta}}$  is also independent of any function of  $\hat{\boldsymbol{\varepsilon}}$ , such as  $\hat{\sigma}^2$ .

PROBLEM 330. *Computer assignment: You run a regression with 3 explanatory variables, no constant term, the sample size is 20, the errors are normally distributed and you know that  $\sigma^2 = 2$ . Plot the density function of  $s^2$ . Hint: The command `dchisq(x, df=25)` returns the density of a Chi-square distribution with 25 degrees of freedom evaluated at  $x$ . But the number 25 was only taken as an example, this is not the number of degrees of freedom you need here.*

• a. *In the same plot, plot the density function of the Theil-Schweitzer estimate. Can one see from the comparison of these density functions why the Theil-Schweitzer estimator has a better MSE?*

ANSWER. Start with the Theil-Schweitzer plot, because it is higher. `> x <- seq(from = 0, to = 6, by = 0.01) > Density <- (19/2)*dchisq((19/2)*x, df=17) > plot(x, Density, type="l", lty=2) > lines(x, (17/2)*dchisq((17/2)*x, df=17)) > title(main = "Unbiased versus Theil-Schw Variance Estimate, 17 d.f.")` □

Now let us derive the maximum likelihood estimator in the case of nonspherical but positive definite covariance matrix. I.e., the model is  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\boldsymbol{\Psi})$ . The density function is

$$(27.0.11) \quad f_{\mathbf{y}}(\mathbf{y}) = (2\pi\sigma^2)^{-n/2} |\det \boldsymbol{\Psi}|^{-1/2} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right).$$

PROBLEM 331. *Derive (27.0.11) as follows: Take a matrix  $\mathbf{P}$  with the property that  $\mathbf{P}\boldsymbol{\varepsilon}$  has covariance matrix  $\sigma^2\mathbf{I}$ . Write down the joint density function of  $\mathbf{P}\boldsymbol{\varepsilon}$ . Since  $\mathbf{y}$  is a linear transformation of  $\boldsymbol{\varepsilon}$ , one can apply the rule for the density function of a transformed random variable.*

ANSWER. Write  $\boldsymbol{\Psi} = \mathbf{Q}\mathbf{Q}^\top$  with  $\mathbf{Q}$  nonsingular and define  $\mathbf{P} = \mathbf{Q}^{-1}$  and  $\mathbf{v} = \mathbf{P}\boldsymbol{\varepsilon}$ . Then  $\mathcal{V}[\mathbf{v}] = \sigma^2\mathbf{P}\mathbf{Q}\mathbf{Q}^\top\mathbf{P}^\top = \sigma^2\mathbf{I}$ , therefore

$$(27.0.12) \quad f_{\mathbf{v}}(\mathbf{v}) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}\mathbf{v}^\top\mathbf{v}\right).$$

For the transformation rule, write  $\mathbf{v}$ , whose density function you know, as a function of  $\mathbf{y}$ , whose density function you want to know.  $\mathbf{v} = \mathbf{P}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ ; therefore the Jacobian matrix is  $\partial\mathbf{v}/\partial\mathbf{y}^\top = \partial(\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{X}\boldsymbol{\beta})/\partial\mathbf{y}^\top = \mathbf{P}$ , or one can see it also element by element

$$(27.0.13) \quad \begin{bmatrix} \frac{\partial v_1}{\partial y_1} & \cdots & \frac{\partial v_1}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial v_n}{\partial y_1} & \cdots & \frac{\partial v_n}{\partial y_n} \end{bmatrix} = \mathbf{P},$$

therefore one has to do two things: first, substitute  $\mathbf{P}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  for  $\mathbf{v}$  in formula (27.0.12), and secondly multiply by the absolute value of the determinant of the Jacobian. Here is how to express the determinant of the Jacobian in terms of  $\boldsymbol{\Psi}$ : From  $\boldsymbol{\Psi}^{-1} = (\mathbf{Q}\mathbf{Q}^\top)^{-1} = (\mathbf{Q}^\top)^{-1}\mathbf{Q}^{-1} = (\mathbf{Q}^{-1})^\top\mathbf{Q}^{-1} = \mathbf{P}^\top\mathbf{P}$  follows  $(\det \mathbf{P})^2 = (\det \boldsymbol{\Psi})^{-1}$ , hence  $|\det \mathbf{P}| = \sqrt{\det \boldsymbol{\Psi}}$ . □

From (27.0.11) one obtains the following log likelihood function: (27.0.14)

$$\log f_{\mathbf{y}}(\mathbf{y}) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2} \ln \det[\boldsymbol{\Psi}] - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

Here, usually not only the elements of  $\boldsymbol{\beta}$  are unknown, but also  $\boldsymbol{\Psi}$  depends on unknown parameters. Instead of concentrating out  $\boldsymbol{\beta}$ , we will first concentrate on  $\sigma^2$ , i.e., we will compute the maximum of this likelihood function over  $\sigma^2$  for a given set of values for the data and the other parameters:

$$(27.0.15) \quad \frac{\partial}{\partial \sigma^2} \log f_{\mathbf{y}}(\mathbf{y}) = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^4}$$

$$(27.0.16) \quad \tilde{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{n}.$$

Whatever the value of  $\boldsymbol{\beta}$  or the values of the unknown parameters in  $\boldsymbol{\Psi}$ ,  $\tilde{\sigma}^2$  is a value of  $\sigma^2$  which, together with the given  $\boldsymbol{\beta}$  and  $\boldsymbol{\Psi}$ , gives the highest value of the likelihood function. If one plugs this  $\tilde{\sigma}^2$  into the likelihood function, one obtains the so-called “concentrated likelihood function” which then only has to be maximized over  $\boldsymbol{\beta}$  and  $\boldsymbol{\Psi}$ :

$$(27.0.17) \quad \log f_{\mathbf{y}}(\mathbf{y}; \tilde{\sigma}^2) = -\frac{n}{2}(1 + \ln 2\pi - \ln n) - \frac{n}{2} \ln(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \frac{1}{2} \ln \det[\boldsymbol{\Psi}]$$

This objective function has to be maximized with respect to  $\boldsymbol{\beta}$  and the parameters entering  $\boldsymbol{\Psi}$ . If  $\boldsymbol{\Psi}$  is known, then this is clearly maximized by the  $\hat{\boldsymbol{\beta}}$  minimizing (19.0.11), therefore the GLS estimator is also the maximum likelihood estimator.

If  $\boldsymbol{\Psi}$  depends on unknown parameters, it is interesting to compare the maximum likelihood estimator with the nonlinear least squares estimator. The objective function minimized by nonlinear least squares is  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ , which is the sum of squares of the innovation parts of the residuals. These two objective functions therefore differ by the factor  $(\det[\boldsymbol{\Psi}])^{\frac{1}{n}}$ , which only matters if there are unknown parameters in  $\boldsymbol{\Psi}$ . Asymptotically, the objective functions are identical.

Using the factorization theorem for sufficient statistics, one also sees easily that  $\hat{\sigma}^2$  and  $\hat{\boldsymbol{\beta}}$  together form sufficient statistics for  $\sigma^2$  and  $\boldsymbol{\beta}$ . For this use the identity

$$(27.0.18) \quad \begin{aligned} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) &= (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \\ &= (n - k)s^2 + (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}). \end{aligned}$$

Therefore the observation  $\mathbf{y}$  enters the likelihood function only through the sufficient statistics  $\hat{\boldsymbol{\beta}}$  and  $s^2$ . The factorization of the likelihood function is therefore the trivial factorization in which that part which does not depend on the unknown parameters but only on the data is unity.

PROBLEM 332. 12 points The log likelihood function in the linear model is given by (27.0.5). Show that the inverse of the information matrix is

$$(27.0.19) \quad \begin{bmatrix} \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} & \mathbf{o} \\ \mathbf{o}^\top & 2\sigma^4/n \end{bmatrix}$$

The information matrix can be obtained in two different ways. Its typical element has the following two forms:

$$(27.0.20) \quad \mathbb{E}\left[\frac{\partial \ln \ell}{\partial \theta_i} \frac{\partial \ln \ell}{\partial \theta_k}\right] = -\mathbb{E}\left[\frac{\partial^2 \ln \ell}{\partial \theta_i \partial \theta_k}\right],$$

or written as matrix derivatives

$$(27.0.21) \quad \mathcal{E}\left[\frac{\partial \ln \ell}{\partial \boldsymbol{\theta}} \frac{\partial \ln \ell}{\partial \boldsymbol{\theta}^\top}\right] = -\mathcal{E}\left[\frac{\partial^2 \ln \ell}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}\right].$$

In our case  $\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\beta} \\ \sigma^2 \end{bmatrix}$ . The expectation is taken under the assumption that the parameter values are the true values. Compute it both ways.

ANSWER. The log likelihood function can be written as

$$(27.0.22) \quad \ln \ell = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2}(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}).$$

The first derivatives were already computed for the maximum likelihood estimators:

$$(27.0.23) \quad \frac{\partial}{\partial \boldsymbol{\beta}^\top} \ln \ell = -\frac{1}{2\sigma^2}(2\mathbf{y}^\top \mathbf{X} + 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}) = \frac{1}{\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{X} = \frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top \mathbf{X}$$

$$(27.0.24) \quad \frac{\partial}{\partial \sigma^2} \ln \ell = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}$$

By the way, one sees that each of these has expected value zero, which is a fact that is needed to prove consistency of the maximum likelihood estimator.

The formula with only one partial derivative will be given first, although it is more tedious:

By doing  $\frac{\partial}{\partial \boldsymbol{\beta}^\top} \left(\frac{\partial}{\partial \boldsymbol{\beta}^\top}\right)^\top$  we get a symmetric  $2 \times 2$  partitioned matrix with the diagonal elements

$$(27.0.25) \quad \mathcal{E}\left[\frac{1}{\sigma^4} \mathbf{X}^\top \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \mathbf{X}\right] = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

and

$$(27.0.26) \quad \mathbb{E}\left[\left(-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}\right)^2\right] = \text{var}\left[-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}\right] = \text{var}\left[\frac{1}{2\sigma^4} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}\right] = \frac{1}{4\sigma^8} 2n\sigma^4 = \frac{n}{2\sigma^4}$$

One of the off-diagonal elements is  $\left(-\frac{n}{2\sigma^4} + \frac{1}{2\sigma^6} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}\right) \boldsymbol{\varepsilon}^\top \mathbf{X}$ . Its expected value is zero:  $\mathcal{E}[\boldsymbol{\varepsilon}] = \mathbf{o}$ , and also  $\mathcal{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top] = \mathbf{o}$  since its  $i$ th component is  $\mathbb{E}[\varepsilon_i \sum_j \varepsilon_j^2] = \sum_j \mathbb{E}[\varepsilon_i \varepsilon_j^2]$ . If  $i \neq j$ , then  $\varepsilon_i$  is independent of  $\varepsilon_j^2$ , therefore  $\mathbb{E}[\varepsilon_i \varepsilon_j^2] = 0 \cdot \sigma^2 = 0$ . If  $i = j$ , we get  $\mathbb{E}[\varepsilon_i^3] = 0$  since  $\varepsilon_i$  has a symmetric distribution.

It is easier if we differentiate once more:

$$(27.0.27) \quad \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \ln \ell = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

$$(27.0.28) \quad \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \sigma^2} \ln \ell = -\frac{1}{\sigma^4} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = -\frac{1}{\sigma^4} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}$$

$$(27.0.29) \quad \frac{\partial^2}{(\partial \sigma^2)^2} \ln \ell = \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}$$

This gives the top matrix in [JHG<sup>+</sup>88, (6.1.24b)]:

$$(27.0.30) \quad \begin{bmatrix} -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} & -\frac{1}{\sigma^4} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}) \\ -\frac{1}{\sigma^4} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta})^\top & \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \end{bmatrix}$$

Now assume that  $\boldsymbol{\beta}$  and  $\sigma^2$  are the true values, take expected values, and reverse the sign. This gives the information matrix

$$(27.0.31) \quad \begin{bmatrix} \sigma^{-2} \mathbf{X}^\top \mathbf{X} & \mathbf{o} \\ \mathbf{o}^\top & n/(2\sigma^4) \end{bmatrix}$$

For the lower righthand side corner we need that  $\mathbb{E}[(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})] = \mathbb{E}[\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}] = n\sigma^2$ .

Taking inverses gives (27.0.19), which is a lower bound for the covariance matrix; we see that  $s^2$  with  $\text{var}[s^2] = 2\sigma^4/(n - k)$  does not attain the bound. However one can show with other means that it is nevertheless efficient.



CHAPTER 28

## Random Regressors

Until now we always assumed that  $\mathbf{X}$  was nonrandom, i.e., the hypothetical repetitions of the experiment used the same  $\mathbf{X}$  matrix. In the nonexperimental sciences, such as economics, this assumption is clearly inappropriate. It is only justified because most results valid for nonrandom regressors can be generalized to the case of random regressors. To indicate that the regressors are random, we will write them as  $\mathbf{X}$ .

### 28.1. Strongest Assumption: Error Term Well Behaved Conditionally on Explanatory Variables

The assumption which we will discuss first is that  $\mathbf{X}$  is random, but the classical assumptions hold *conditionally on  $\mathbf{X}$* , i.e., the conditional expectation  $\mathcal{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{o}$ , and the conditional variance-covariance matrix  $\mathcal{V}[\boldsymbol{\varepsilon}|\mathbf{X}] = \sigma^2\mathbf{I}$ . In this situation, the least squares estimator has all the classical properties *conditionally on  $\mathbf{X}$* , for instance  $\mathcal{E}[\hat{\boldsymbol{\beta}}|\mathbf{X}] = \boldsymbol{\beta}$ ,  $\mathcal{V}[\hat{\boldsymbol{\beta}}|\mathbf{X}] = \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}$ ,  $\mathbb{E}[s^2|\mathbf{X}] = \sigma^2$ , etc.

Moreover, certain properties of the Least Squares estimator remain valid *unconditionally*. An application of the law of iterated expectations shows that the least squares estimator  $\hat{\boldsymbol{\beta}}$  is still unbiased. Start with (18.0.7):

$$(28.1.1) \quad \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\varepsilon}$$

$$(28.1.2) \quad \mathcal{E}[\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}|\mathbf{X}] = \mathcal{E}[(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{\varepsilon}|\mathbf{X}] = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathcal{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{o}.$$

$$(28.1.3) \quad \mathcal{E}[\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}] = \mathcal{E}[\mathcal{E}[\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}|\mathbf{X}]] = \mathbf{o}.$$

PROBLEM 333. 1 point In the model with random explanatory variables  $\mathbf{X}$  you are considering an estimator  $\tilde{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$ . Which statement is stronger:  $\mathcal{E}[\tilde{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ , or  $\mathcal{E}[\tilde{\boldsymbol{\beta}}|\mathbf{X}] = \boldsymbol{\beta}$ . Justify your answer.

ANSWER. The second statement is stronger. The first statement follows from the second by the law of iterated expectations.  $\square$

PROBLEM 334. 2 points Assume the regressors  $\mathbf{X}$  are random, and the classical assumptions hold conditionally on  $\mathbf{X}$ , i.e.,  $\mathcal{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{o}$  and  $\mathcal{V}[\boldsymbol{\varepsilon}|\mathbf{X}] = \sigma^2\mathbf{I}$ . Show that  $s^2$  is an unbiased estimate of  $\sigma^2$ .

ANSWER. From the theory with nonrandom explanatory variables follows  $\mathbb{E}[s^2|\mathbf{X}] = \sigma^2$ . Therefore  $\mathbb{E}[s^2] = \mathbb{E}[\mathbb{E}[s^2|\mathbf{X}]] = \mathbb{E}[\sigma^2] = \sigma^2$ . In words: if the expectation conditional on  $\mathbf{X}$  does not depend on  $\mathbf{X}$ , then it is also the unconditional expectation.

The law of iterated expectations can also be used to compute the unconditional  $MSE$  matrix of  $\hat{\boldsymbol{\beta}}$ :

$$(28.1.4) \quad MSE[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}] = \mathcal{E}[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top]$$

$$(28.1.5) \quad = \mathcal{E}[\mathcal{E}[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top|\mathbf{X}]]$$

$$(28.1.6) \quad = \mathcal{E}[\sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}]$$

$$(28.1.7) \quad = \sigma^2\mathcal{E}[(\mathbf{X}^\top\mathbf{X})^{-1}].$$

PROBLEM 335. 2 points Show that  $s^2(\mathbf{X}^\top\mathbf{X})^{-1}$  is unbiased estimator of  $MSE[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}]$ .

ANSWER.

$$(28.1.8) \quad \mathcal{E}[s^2(\mathbf{X}^\top\mathbf{X})^{-1}] = \mathcal{E}[\mathcal{E}[s^2(\mathbf{X}^\top\mathbf{X})^{-1}|\mathbf{X}]]$$

$$(28.1.9) \quad = \mathcal{E}[\sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}]$$

$$(28.1.10) \quad = \sigma^2\mathcal{E}[(\mathbf{X}^\top\mathbf{X})^{-1}]$$

$$(28.1.11) \quad = MSE[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}] \quad \text{by (28.1.7).}$$

The Gauss-Markov theorem generalizes in the following way: Say  $\tilde{\boldsymbol{\beta}}$  is an estimator, linear in  $\mathbf{y}$ , but not necessarily in  $\mathbf{X}$ , satisfying  $\mathcal{E}[\tilde{\boldsymbol{\beta}}|\mathbf{X}] = \boldsymbol{\beta}$  (which is stronger than unbiasedness); then  $MSE[\tilde{\boldsymbol{\beta}}; \boldsymbol{\beta}] \geq MSE[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}]$ . Proof is immediate: we know by the usual Gauss-Markov theorem that  $MSE[\tilde{\boldsymbol{\beta}}; \boldsymbol{\beta}|\mathbf{X}] \geq MSE[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}|\mathbf{X}]$ , and taking the expected values will preserve this inequality:  $\mathcal{E}[MSE[\tilde{\boldsymbol{\beta}}; \boldsymbol{\beta}|\mathbf{X}]] \geq \mathcal{E}[MSE[\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}|\mathbf{X}]]$ , but this expected value is exactly the unconditional  $MSE$ .

The assumption  $\mathcal{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{o}$  can also be written  $\mathcal{E}[\mathbf{y}|\mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$ , and  $\mathcal{V}[\boldsymbol{\varepsilon}|\mathbf{X}] = \sigma^2\mathbf{I}$  can also be written as  $\mathcal{V}[\mathbf{y}|\mathbf{X}] = \sigma^2\mathbf{I}$ . Both of these are assumptions about the conditional distribution  $\mathbf{y}|\mathbf{X} = \mathbf{X}$  for all  $\mathbf{X}$ . This suggests the following broadening of the regression paradigm:  $\mathbf{y}$  and  $\mathbf{X}$  are jointly distributed random variables, and one is interested how  $\mathbf{y}|\mathbf{X} = \mathbf{X}$  depends on  $\mathbf{X}$ . If the expected value of this distribution depends linearly, and the variance of this distribution is constant, then this is the linear regression model discussed above. But the expected value might also depend on  $\mathbf{X}$  in a nonlinear fashion (nonlinear least squares), and the variance may not be constant—in which case the intuition that  $\mathbf{y}$  is some function of  $\mathbf{X}$  plus some error term may no longer be appropriate;  $\mathbf{y}$  may for instance be the outcome of a binomial choice, the probability of which depends on  $\mathbf{X}$  (see chapter ??; the generalized linear model).

## 28.2. Contemporaneously Uncorrelated Disturbances

In many situations with random regressors, the condition  $\mathcal{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{o}$  is not satisfied. Instead, the columns of  $\mathbf{X}$  are contemporaneously uncorrelated with  $\boldsymbol{\varepsilon}$ , but they may be correlated with past values of  $\boldsymbol{\varepsilon}$ . The main example here is regression with a lagged dependent variable. In this case, OLS is no longer unbiased, but asymptotically it still has all the good properties, it is asymptotically normal with the covariance matrix which one would expect. Asymptotically, the computer printout is still valid. This is a very important result, which is often used in econometrics, but most econometrics textbooks do not even start to prove it. There is a proof in [Kme86, pp. 749–757], and one in [Mal80, pp. 535–539].

**PROBLEM 336.** *Since least squares with random regressors is appropriate whenever the disturbances are contemporaneously uncorrelated with the explanatory variables, a friend of yours proposes to test for random explanatory variables by checking whether the sample correlation coefficients between the residuals and the explanatory variables is significantly different from zero or not. Is this an appropriate statistic?*

**ANSWER.** No. The sample correlation coefficients are always zero!

□

## 28.3. Disturbances Correlated with Regressors in Same Observation

But if  $\boldsymbol{\varepsilon}$  is contemporaneously correlated with  $\mathbf{X}$ , then OLS is inconsistent. This can be the case in some dynamic processes (lagged dependent variable as regressor, and autocorrelated errors, see question ??), when there are, in addition to the relation which one wants to test with the regression, other relations making the righthand side variables dependent on the lefthand side variable, or when the righthand side variables are measured with errors. This is usually the case in economics, and econometrics has developed the technique of simultaneous equations estimation to deal with it.

**PROBLEM 337.** *3 points What does one have to watch out for if some of the regressors are random?*

## CHAPTER 29

## The Mahalanobis Distance

Everything in this chapter is unpublished work, presently still in draft form. The aim is to give a motivation for the least squares objective function in terms of an initial measure of precision. The case of prediction is mathematically simpler than that of estimation, therefore this chapter will only discuss prediction. We assume that the joint distribution of  $\mathbf{y}$  and  $\mathbf{z}$  has the form

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} \sim \begin{bmatrix} \mathbf{X} \\ \mathbf{W} \end{bmatrix} \boldsymbol{\beta}, \sigma^2 \begin{bmatrix} \boldsymbol{\Omega}_{\mathbf{y}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{y}\mathbf{z}} \\ \boldsymbol{\Omega}_{\mathbf{z}\mathbf{y}} & \boldsymbol{\Omega}_{\mathbf{z}\mathbf{z}} \end{bmatrix}, \quad (29.0.1)$$

$\sigma^2 > 0$ , otherwise unknown  
 $\boldsymbol{\beta}$  unknown as well.

$\mathbf{y}$  is observed but  $\mathbf{z}$  is not and has to be predicted. But assume we are not interested in the  $\mathcal{MSE}$  since we do the experiment only once. We want to predict  $\mathbf{z}$  in such a way that, whatever the true value of  $\boldsymbol{\beta}$ , the predicted value  $\mathbf{z}^*$  “blends in” best with the given data  $\mathbf{y}$ .

There is an important conceptual difference between this criterion and the one based on the  $\mathcal{MSE}$ . The present criterion cannot be applied until after the data are known, therefore it is called a “final” criterion as opposed to the “initial” criterion of the  $\mathcal{MSE}$ . See Barnett [Bar82, pp. 157–159] for a good discussion of these issues.

How do we measure the degree to which a given data set “blend in,” i.e., are not outliers for a given distribution? Hypothesis testing uses this criterion. The most often-used testing principle is: reject the null hypothesis if the observed value of a certain statistic is too much an outlier for the distribution which this statistic would have under the null hypothesis. If the statistic is a scalar, and if under the null hypothesis this statistic has expected value  $\mu$  and standard deviation  $\sigma$ , then one often uses an estimate of  $|x - \mu|/\sigma$ , the number of standard deviations the observed value is away from the mean, to measure the “distance” of the observed value  $x$  from the distribution  $(\mu, \sigma^2)$ . The Mahalanobis distance generalizes this concept to the case that the test statistic is a vector random variable.

## 29.1. Definition of the Mahalanobis Distance

Since it is mathematically more convenient to work with the *squared* distance than with the distance itself, we will make the following thought experiment to

motivate the Mahalanobis distance. How could one generalize the squared scalar distance  $(y - \mu)^2/\sigma^2$  for the distance of a vector value  $\mathbf{y}$  from the distribution of the vector random variable  $\mathbf{y} \sim (\boldsymbol{\mu}, \sigma^2\boldsymbol{\Omega})$ ? If all  $y_i$  have same variance  $\sigma^2$ , i.e.  $\boldsymbol{\Omega} = \mathbf{I}$ , one might measure the squared distance of  $\mathbf{y}$  from the distribution  $(\boldsymbol{\mu}, \sigma^2)$  by  $\frac{1}{\sigma^2} \max_i (y_i - \mu_i)^2$ , but since the maximum from two trials is bigger than the value from one trial only, one should divide this perhaps by the expected value of such a maximum. If the variances are different, say  $\sigma_i^2$ , one might want to look at the number of standard deviations which the “worst” component of  $\mathbf{y}$  is away from what would be its mean if  $\mathbf{y}$  were an observation of  $\mathbf{y}$ , i.e., the squared distance of the observed vector from the distribution would be  $\max_i \frac{(y_i - \mu_i)^2}{\sigma_i^2}$ , again normalized by the expected value.

The principle actually used by the Mahalanobis distance goes only a small step further than the examples just cited. It is coordinate-free, i.e., any linear combinations of the elements of  $\mathbf{y}$  are considered on equal footing with these elements themselves. In other words, it does not distinguish between variates and variables. The distance of a given vector value from a certain multivariate distribution is defined to be the distance of the “worst” *linear combination* of the elements of this vector from the univariate distribution of this linear combination, normalized in such a way that the expected value of this distance is 1.

DEFINITION 29.1.1. Given a random  $n$ -vector  $\mathbf{y}$  which has expected value  $\boldsymbol{\mu}$  and nonsingular covariance matrix. The squared “Mahalanobis distance” or “statistical distance” of the observed value  $\mathbf{y}$  from the distribution of  $\mathbf{y}$  is defined to be

$$(29.1.1) \quad \text{MHD}[\mathbf{y}; \mathbf{y}] = \frac{1}{n} \max_{\mathbf{g}} \frac{(\mathbf{g}^\top \mathbf{y} - \mathbf{E}[\mathbf{g}^\top \mathbf{y}])^2}{\text{var}[\mathbf{g}^\top \mathbf{y}]}.$$

If the denominator  $\text{var}[\mathbf{g}^\top \mathbf{y}]$  is zero, then  $\mathbf{g} = \mathbf{o}$ , therefore the numerator is zero as well. In this case the fraction is defined to be zero.

THEOREM 29.1.2. Let  $\mathbf{y}$  be a vector random variable with  $\mathcal{E}[\mathbf{y}] = \boldsymbol{\mu}$  and  $\mathcal{V}[\mathbf{y}] = \sigma^2\boldsymbol{\Omega}$ ,  $\sigma^2 > 0$  and  $\boldsymbol{\Omega}$  positive definite. The squared Mahalanobis distance of the value  $\mathbf{y}$  from the distribution of  $\mathbf{y}$  is equal to

$$(29.1.2) \quad \text{MHD}[\mathbf{y}; \mathbf{y}] = \frac{1}{n\sigma^2} (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

PROOF. (29.1.2) is a simple consequence of (25.4.4). It is also somewhat intuitive since the righthand side of (29.1.2) can be considered a division of the square of  $\mathbf{y} - \boldsymbol{\mu}$  by the covariance matrix of  $\mathbf{y}$ .

The Mahalanobis distance is an asymmetric measure; a large value indicates a bad fit of the hypothetical population to the observation, while a value of, say, 1 does not necessarily indicate a better fit than a value of 1.

PROBLEM 338. Let  $\mathbf{y}$  be a random  $n$ -vector with expected value  $\boldsymbol{\mu}$  and nonsingular covariance matrix  $\sigma^2\boldsymbol{\Omega}$ . Show that the expected value of the Mahalanobis distance of the observations of  $\mathbf{y}$  from the distribution of  $\mathbf{y}$  is 1, i.e.,

$$(29.1.3) \quad E[\text{MHD}[\mathbf{y}; \mathbf{y}]] = 1$$

ANSWER.

(29.1.4)

$$E\left[\frac{1}{n\sigma^2}(\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Omega}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right] = E\left[\text{tr}\left(\frac{1}{n\sigma^2}\boldsymbol{\Omega}^{-1}(\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})^\top\right)\right] = \frac{1}{n}\text{tr}(\mathbf{I}) = 1.$$

□

(29.1.2) is, up to a constant factor, the quadratic form in the exponent of the normal density function of  $\mathbf{y}$ . For a normally distributed  $\mathbf{y}$ , therefore, all observations located on the same density contour have equal distance from the distribution.

The Mahalanobis distance is also defined if the covariance matrix of  $\mathbf{y}$  is singular. In this case, certain nonzero linear combinations of the elements of  $\mathbf{y}$  are known with certainty. Certain vectors can therefore not possibly be realizations of  $\mathbf{y}$ , i.e., the set of realizations of  $\mathbf{y}$  does not fill the whole  $\mathbb{R}^n$ .

PROBLEM 339. 2 points The random vector  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$  has mean  $\begin{bmatrix} 1 \\ 2 \\ -3 \end{bmatrix}$  and covariance matrix  $\frac{1}{3} \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix}$ . Is this covariance matrix singular? If so, give a linear combination of the elements of  $\mathbf{y}$  which is known with certainty. And give a value which can never be a realization of  $\mathbf{y}$ . Prove everything you state.

ANSWER. Yes, it is singular;

$$(29.1.5) \quad \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

I.e.,  $y_1 + y_2 + y_3 = 0$  because its variance is 0 and its mean is zero as well since  $[1 \ 1 \ 1] \begin{bmatrix} 1 \\ 2 \\ -3 \end{bmatrix} = 0$ . □

DEFINITION 29.1.3. Given a vector random variable  $\mathbf{y}$  which has a mean and a covariance matrix. A value  $\mathbf{y}$  has infinite statistical distance from this random variable, i.e., it cannot possibly be a realization of this random variable, if a vector of coefficients  $\mathbf{g}$  exists such that  $\text{var}[\mathbf{g}^\top \mathbf{y}] = 0$  but  $\mathbf{g}^\top \mathbf{y} \neq \mathbf{g}^\top E[\mathbf{y}]$ . If such a  $\mathbf{g}$  does not exist, then the squared Mahalanobis distance of  $\mathbf{y}$  from  $\mathbf{y}$  is defined as in (29.1.1), with  $n$  replaced by  $\text{rank}[\boldsymbol{\Omega}]$ . If the denominator in (29.1.1) is zero, then it no longer necessarily follows that  $\mathbf{g} = \mathbf{o}$  but it nevertheless follows that the numerator is zero, and the fraction should in this case again be considered zero.

If  $\boldsymbol{\Omega}$  is singular, then the inverse  $\boldsymbol{\Omega}^{-1}$  in formula (29.1.2) must be replaced by a “g-inverse.” A g-inverse of a matrix  $\mathbf{A}$  is any matrix  $\mathbf{A}^-$  which satisfies  $\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}$ . G-inverses always exist, but they are usually not unique.

PROBLEM 340.  $a$  is a scalar. What is its g-inverse  $a^-$ ?

THEOREM 29.1.4. Let  $\mathbf{y}$  be a random variable with  $E[\mathbf{y}] = \boldsymbol{\mu}$  and  $\mathcal{V}[\mathbf{y}] = \sigma^2$ ,  $\sigma^2 > 0$ . If it is not possible to express the vector  $\mathbf{y}$  in the form  $\mathbf{y} = \boldsymbol{\mu} + \boldsymbol{\Omega}\mathbf{a}$  for some  $\mathbf{a}$ , then the squared Mahalanobis distance of  $\mathbf{y}$  from the distribution of  $\mathbf{y}$  is infinite, i.e.,  $\text{MHD}[\mathbf{y}; \mathbf{y}] = \infty$ ; otherwise

$$(29.1.6) \quad \text{MHD}[\mathbf{y}; \mathbf{y}] = \frac{1}{\sigma^2 \text{rank}[\boldsymbol{\Omega}]}(\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Omega}^-(\mathbf{y} - \boldsymbol{\mu})$$

Now we will discuss how a given observation vector can be extended by additional observations in such a way that the Mahalanobis distance of the whole vector from its distribution is minimized.

## CHAPTER 30

## Interval Estimation

We will first show how the least squares principle can be used to construct confidence regions, and then we will derive the properties of these confidence regions.

## 30.1. A Basic Construction Principle for Confidence Regions

The least squares objective function, whose minimum argument gave us the BLUE, naturally allows us to generate *confidence intervals* or higher-dimensional *confidence regions*. A confidence region for  $\beta$  based on  $\mathbf{y} = \mathbf{X}\beta + \varepsilon$  can be constructed as follows:

- Draw the OLS estimate  $\hat{\beta}$  into  $k$ -dimensional space; it is the vector which minimizes  $SSE = (\mathbf{y} - \mathbf{X}\hat{\beta})^\top(\mathbf{y} - \mathbf{X}\hat{\beta})$ .
- For every other vector  $\tilde{\beta}$  one can define the sum of squared errors associated with that vector as  $SSE_{\tilde{\beta}} = (\mathbf{y} - \mathbf{X}\tilde{\beta})^\top(\mathbf{y} - \mathbf{X}\tilde{\beta})$ . Draw the level hypersurfaces (if  $k = 2$ : level lines) of this function. These are ellipsoids centered on  $\hat{\beta}$ .
- Each of these ellipsoids is a confidence region for  $\beta$ . Different confidence regions differ by their coverage probabilities.
- If one is only interested in certain coordinates of  $\beta$  and not in the others, or in some other linear transformation  $\beta$ , then the corresponding confidence regions are the corresponding transformations of this ellipse. Geometrically this can best be seen if this transformation is an orthogonal projection; then the confidence ellipse of the transformed vector  $\mathbf{R}\beta$  is also a projection or “shadow” of the confidence region for the whole vector. Projections of the same confidence region have the same confidence level, independent of the direction in which this projection goes.

The confidence regions for  $\beta$  with coverage probability  $\pi$  will be written here as  $B_{\beta;\pi}$  or, if we want to make its dependence on the observation vector  $\mathbf{y}$  explicit,  $B_{\beta;\pi}(\mathbf{y})$ . These confidence regions are level lines of the SSE, and mathematically, it is advantageous to define these level lines by their level relative to the minimum level, i.e., as the set of all  $\tilde{\beta}$  for which the quotient of the attained  $SSE_{\tilde{\beta}} =$

$(\mathbf{y} - \mathbf{X}\tilde{\beta})^\top(\mathbf{y} - \mathbf{X}\tilde{\beta})$  divided by the smallest possible  $SSE = (\mathbf{y} - \mathbf{X}\hat{\beta})^\top(\mathbf{y} - \mathbf{X}\hat{\beta})$  is smaller or equal a given number. In formulas,

$$(30.1.1) \quad \tilde{\beta} \in B_{\beta;\pi}(\mathbf{y}) \iff \frac{(\mathbf{y} - \mathbf{X}\tilde{\beta})^\top(\mathbf{y} - \mathbf{X}\tilde{\beta})}{(\mathbf{y} - \mathbf{X}\hat{\beta})^\top(\mathbf{y} - \mathbf{X}\hat{\beta})} \leq c_{\pi;n-k,k}$$

It will be shown below, in the discussion following (30.2.1), that  $c_{\pi;n-k,k}$  only depends on  $\pi$  (the confidence level),  $n - k$  (the degrees of freedom in the regression), and  $k$  (the dimension of the confidence region).

To get a geometric intuition of this principle, look at the case  $k = 2$ , in which the parameter vector  $\beta$  has only two components. For each possible value  $\tilde{\beta}$  of the parameter vector, the associated sum of squared errors is  $SSE_{\tilde{\beta}} = (\mathbf{y} - \mathbf{X}\tilde{\beta})^\top(\mathbf{y} - \mathbf{X}\tilde{\beta})$ . This is a quadratic function of  $\tilde{\beta}$ , whose level lines form concentric ellipses shown in Figure 1. The center of these ellipses is the unconstrained least squares estimate. Each of the ellipses is a confidence region for  $\beta$  for a different confidence level.

If one needs a confidence region not for the whole vector  $\beta$  but, say, for  $i$  linearly independent linear combinations  $\mathbf{R}\beta$  (here  $\mathbf{R}$  is a  $i \times k$  matrix with full row rank), then the above principle applies in the following way: the vector  $\tilde{\mathbf{u}}$  lies in the confidence region for  $\mathbf{R}\beta$  generated by  $\mathbf{y}$  for confidence level  $\pi$ , notation  $B_{\mathbf{R}\beta;\pi}$ , if and only if there is a  $\tilde{\beta}$  in the confidence region (30.1.1) (with the parameters adjusted to reflect the dimensionality of  $\tilde{\mathbf{u}}$ ) which satisfies  $\mathbf{R}\tilde{\beta} = \tilde{\mathbf{u}}$ :

$$(30.1.2) \quad \tilde{\mathbf{u}} \in B_{\mathbf{R}\beta;\pi}(\mathbf{y}) \iff \text{exist } \tilde{\beta} \text{ with } \tilde{\mathbf{u}} = \mathbf{R}\tilde{\beta} \text{ and } \frac{(\mathbf{y} - \mathbf{X}\tilde{\beta})^\top(\mathbf{y} - \mathbf{X}\tilde{\beta})}{(\mathbf{y} - \mathbf{X}\hat{\beta})^\top(\mathbf{y} - \mathbf{X}\hat{\beta})} \leq c_{\pi;n-k,k}$$

**PROBLEM 341.** *Why does one have to change the value of  $c$  when one goes over to the projections of the confidence regions?*

**ANSWER.** Because the projection is a many-to-one mapping, and vectors which are not in the original ellipsoid may still end up in the projection.

Again let us illustrate this with the 2-dimensional case in which the confidence region for  $\beta$  is an ellipse, as drawn in Figure 1, called  $B_{\beta;\pi}(\mathbf{y})$ . Starting with the ellipse, the above criterion defines individual confidence intervals for linear combinations  $u = \mathbf{r}^\top\beta$  by the rule:  $\tilde{u} \in B_{\mathbf{r}^\top\beta;\pi}(\mathbf{y})$  iff a  $\tilde{\beta} \in B_{\beta;\pi}(\mathbf{y})$  exists with  $\mathbf{r}^\top\tilde{\beta} = \tilde{u}$ . For  $\mathbf{r} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ , this interval is simply the projection of the ellipse on the horizontal axis, and for  $\mathbf{r} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$  it is the projection on the vertical axis.

The same argument applies for all vectors  $\mathbf{r}$  with  $\mathbf{r}^\top\mathbf{r} = 1$ . The inner product of two vectors is the length of the first vector times the length of the projection of the second vector on the first. If  $\mathbf{r}^\top\mathbf{r} = 1$ , therefore,  $\mathbf{r}^\top\tilde{\beta}$  is simply the length of the orthogonal projection of  $\tilde{\beta}$  on the line generated by the vector  $\mathbf{r}$ . Therefore

the confidence interval for  $\mathbf{r}^\top \boldsymbol{\beta}$  is simply the projection of the ellipse on the line generated by  $\mathbf{r}$ . (This projection is sometimes called the “shadow” of the ellipse.)

The confidence region for  $\mathbf{R}\boldsymbol{\beta}$  can also be defined as follows:  $\tilde{\mathbf{u}}$  lies in this confidence region if and only if the “best”  $\hat{\boldsymbol{\beta}}$  which satisfies  $\mathbf{R}\hat{\boldsymbol{\beta}} = \tilde{\mathbf{u}}$  lies in the confidence region (30.1.1), this best  $\hat{\boldsymbol{\beta}}$  being, of course, the constrained least squares estimate subject to the constraint  $\mathbf{R}\boldsymbol{\beta} = \tilde{\mathbf{u}}$ , whose formula is given by (22.3.13). The confidence region for  $\mathbf{R}\boldsymbol{\beta}$  consists therefore of all  $\tilde{\mathbf{u}}$  for which the constrained least squares estimate  $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \tilde{\mathbf{u}})$  satisfies condition (30.1.1):

$$(30.1.3) \quad \tilde{\mathbf{u}} \in B_{\mathbf{R}\boldsymbol{\beta}}(\mathbf{y}) \iff \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})} \leq c_{\pi;n-k,i}$$

One can also write it as

$$(30.1.4) \quad \tilde{\mathbf{u}} \in B_{\mathbf{R}\boldsymbol{\beta}}(\mathbf{y}) \iff \frac{\text{SSE}_{\text{constrained}}}{\text{SSE}_{\text{unconstrained}}} \leq c_{\pi;n-k,i}$$

i.e., those  $\tilde{\mathbf{u}}$  are in the confidence region which, if imposed as a constraint on the regression, will not make the SSE too much bigger.

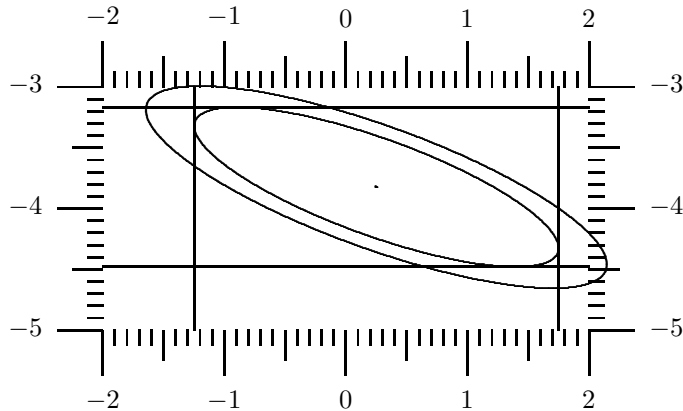


FIGURE 1. Confidence Ellipse with “Shadows”

In order to transform (30.1.3) into a mathematically more convenient form, write it as

$$\tilde{\mathbf{u}} \in B_{\mathbf{R}\boldsymbol{\beta};\pi}(\mathbf{y}) \iff \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) - (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})} \leq c_{\pi;n-k,i} - 1$$

and then use (22.7.2) to get

$$(30.1.5) \quad \tilde{\mathbf{u}} \in B_{\mathbf{R}\boldsymbol{\beta};\pi}(\mathbf{y}) \iff \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \tilde{\mathbf{u}})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \tilde{\mathbf{u}})}{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})} \leq c_{\pi;n-k,i} - 1$$

This formula has the great advantage that  $\hat{\boldsymbol{\beta}}$  no longer appears in it. The condition whether  $\tilde{\mathbf{u}}$  belongs to the confidence region is here formulated in terms of  $\hat{\boldsymbol{\beta}}$  alone.

PROBLEM 342. Using (14.2.12), show that (30.1.1) can be rewritten as

$$(30.1.6) \quad \tilde{\boldsymbol{\beta}} \in B_{\boldsymbol{\beta};\pi}(\mathbf{y}) \iff \frac{(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})}{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})} \leq c_{\pi;n-k,k} - 1$$

Verify that this is the same as (30.1.5) in the special case  $\mathbf{R} = \mathbf{I}$ .

PROBLEM 343. You have run a regression with intercept, but you are not interested in the intercept per se but need a joint confidence region for all slope parameters. Using the notation of Problem 304, show that this confidence region has the form

$$(30.1.7) \quad \tilde{\boldsymbol{\beta}} \in B_{\boldsymbol{\beta};\pi}(\mathbf{y}) \iff \frac{(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})}{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})} \leq c_{\pi;n-k,k-1} - 1$$

I.e., we are sweeping the means out of both regressors and dependent variables, and then we act as if the regression never had an intercept and use the formula for full parameter vector (30.1.6) for these transformed data (except that the number of degrees of freedom  $n-k$  still reflects the intercept as one of the explanatory variables).

ANSWER. Write the full parameter vector as  $\begin{bmatrix} \alpha \\ \boldsymbol{\beta} \end{bmatrix}$  and  $\mathbf{R} = \begin{bmatrix} \mathbf{o} & \mathbf{I} \end{bmatrix}$ . Use (30.1.5) but instead of  $\tilde{\mathbf{u}}$  write  $\tilde{\boldsymbol{\beta}}$ . The only tricky part is the following which uses (23.0.37):

$$(30.1.8) \quad \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top = \begin{bmatrix} \mathbf{o} & \mathbf{I} \end{bmatrix} \begin{bmatrix} 1/n + \bar{\mathbf{x}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \bar{\mathbf{x}} & -\bar{\mathbf{x}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \\ -(\mathbf{X}^\top \mathbf{X})^{-1} \bar{\mathbf{x}} & (\mathbf{X}^\top \mathbf{X})^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{o}^\top \\ \mathbf{I} \end{bmatrix} = (\mathbf{X}^\top \mathbf{X})^{-1}$$

The denominator is  $(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ , but since  $\hat{\boldsymbol{\beta}} = \bar{\mathbf{y}} - \bar{\mathbf{x}}^\top \hat{\boldsymbol{\beta}}$ , see problem 204, the denominator can be rewritten as  $(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ .

PROBLEM 344. 3 points We are in the simple regression  $\mathbf{y}_t = \alpha + \beta x_t + \epsilon_t$ . One draws, for every value of  $x$ , a 95% confidence interval for  $\alpha + \beta x$ , one gets a “confidence band” around the fitted line, as shown in Figure 2. Is the probability that this confidence band covers the true regression line over its whole length equal to 95%, greater than 95%, or smaller than 95%? Give a good verbal reasoning for your answer. You should make sure that your explanation is consistent with the fact that the confidence interval is random and the true regression line is fixed.

FIGURE 2. Confidence Band for Regression Line

### 30.2. Coverage Probability of the Confidence Regions

The probability that any *given known value*  $\tilde{\mathbf{u}}$  lies in the confidence region (30.1.3) depends on the unknown  $\beta$ . But we will show now that the “coverage probability” of the region, i.e., the probability with which the confidence region contains the *unknown true value*  $\mathbf{u} = \mathbf{R}\beta$ , does *not* depend on any unknown parameters.

To get the coverage probability, we must substitute  $\tilde{\mathbf{u}} = \mathbf{R}\beta$  (where  $\beta$  is the true parameter value) in (30.1.5). This gives

$$(30.2.1) \quad \mathbf{R}\beta \in B_{\mathbf{R}\beta;\pi}(\mathbf{y}) \iff \frac{(\mathbf{R}\hat{\beta} - \mathbf{R}\beta)^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\beta} - \mathbf{R}\beta)}{(\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta})} \leq c_{\pi;n-k,i} - 1$$

Let us look at numerator and denominator separately. Under the Normality assumption,  $\mathbf{R}\hat{\beta} \sim N(\mathbf{R}\beta, \sigma^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)$ . Therefore, by (7.4.9), the distribution of the numerator of (30.2.1) is

$$(30.2.2) \quad (\mathbf{R}\hat{\beta} - \mathbf{R}\beta)^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\beta} - \mathbf{R}\beta) \sim \sigma^2 \chi_i^2.$$

This probability distribution only depends on one unknown parameter, namely,  $\sigma^2$ . Regarding the denominator, remember that, by (18.4.2),  $(\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) = \boldsymbol{\varepsilon}^\top \mathbf{M} \boldsymbol{\varepsilon}$ , and if we apply (7.4.9) to this we can see that

$$(30.2.3) \quad (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) \sim \sigma^2 \chi_{n-k}^2$$

Furthermore, numerator and denominator are independent. To see this, look first at  $\hat{\beta}$  and  $\hat{\varepsilon}$ . By Problem 246 they are uncorrelated, and since they are also jointly Normal, it follows that they are independent. If  $\hat{\beta}$  and  $\hat{\varepsilon}$  are independent, any functions of  $\hat{\beta}$  are independent of any functions of  $\hat{\varepsilon}$ . The numerator in the test

statistic (30.2.1) is a function of  $\hat{\beta}$  and the denominator is a function of  $\hat{\varepsilon}$ ; therefore they are independent, as claimed. Lastly, if we divide numerator by denominator the unknown “nuisance parameter”  $\sigma^2$  in their probability distributions cancels out, i.e., the distribution of the quotient is fully known.

To sum up: if  $\tilde{\mathbf{u}}$  is the true value  $\tilde{\mathbf{u}} = \mathbf{R}\beta$ , then the test statistic in (30.2.1) can no longer be observed, but its *distribution* is known; it is a  $\chi_i^2$  divided by independent  $\chi_{n-k}^2$ . Therefore, for every value  $c$ , the probability that the confidence region (30.1.5) contains the true  $\mathbf{R}\beta$  can be computed, and conversely, for any desired coverage probability, the appropriate critical value  $c$  can be computed. As claimed, this critical value only depends on the confidence level  $\pi$  and  $n - k$  and  $i$ .

### 30.3. Conventional Formulas for the Test Statistics

In order to get this test statistic into the form in which it is conventionally tabulated, we must divide both numerator and denominator of (30.1.5) by their degrees of freedom, to get a  $\chi_i^2/i$  divided by an independent  $\chi_{n-k}^2/(n - k)$ . This quotient is called a F-distribution with  $i$  and  $n - k$  degrees of freedom.

The F-distribution is defined as  $F_{i,j} = \frac{\chi_i^2/i}{\chi_j^2/j}$  instead of the seemingly simpler formula  $\frac{\chi_i^2}{\chi_j^2}$ , because the division by the degrees of freedom makes all F-distributions and the associated critical values similar; an observed value below 4 is insignificant but greater values may be significant depending on the number of parameters.

Therefore, instead of (30.2.1), the condition deciding whether a given vector  $\tilde{\mathbf{u}}$  lies in the confidence region for  $\mathbf{R}\beta$  with confidence level  $\pi = 1 - \alpha$  is formulated as follows

$$(30.3.1) \quad \frac{(\text{SSE}_{\text{constrained}} - \text{SSE}_{\text{unconstrained}})/\text{number of constraints}}{\text{SSE}_{\text{unconstr.}}/(\text{numb. of obs.} - \text{numb. of coeff. in unconstr. model})} \leq F_{(i,n-k;\alpha)}$$

Here the constrained SSE is the SSE in the model estimated with the constraint  $\mathbf{R}\beta = \tilde{\mathbf{u}}$  imposed, and  $F_{(i,n-k;\alpha)}$  is the upper  $\alpha$  quantile of the F distribution with  $i$  and  $n - k$  degrees of freedom, i.e., it is that scalar  $c$  for which a random variable which has a F distribution with  $i$  and  $n - k$  degrees of freedom satisfies  $\Pr[F \geq c] = \alpha$ .

### 30.4. Interpretation in terms of Studentized Mahalanobis Distance

The division of numerator and denominator by their degrees of freedom also gives us a second intuitive interpretation of the test statistic in terms of the Mahalanobis distance, see chapter 29. If one divides the denominator by its degrees of freedom one gets an unbiased estimate of  $\sigma^2$

$$(30.4.1) \quad s^2 = \frac{1}{n - k} (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}).$$

Therefore from (30.1.5) one gets the following alternative formula for the joint confidence region  $B(\mathbf{y})$  for the vector parameter  $\mathbf{u} = \mathbf{R}\boldsymbol{\beta}$  for confidence level  $\pi = 1 - \alpha$ :

$$(30.4.2) \quad \tilde{\mathbf{u}} \in B_{\mathbf{R}\boldsymbol{\beta}; 1-\alpha}(\mathbf{y}) \iff \frac{1}{s^2}(\mathbf{R}\hat{\boldsymbol{\beta}} - \tilde{\mathbf{u}})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \tilde{\mathbf{u}}) \leq iF_{(i, n-k; \alpha)}$$

Here  $\hat{\boldsymbol{\beta}}$  is the least squares estimator of  $\boldsymbol{\beta}$ , and  $s^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / (n - k)$  the unbiased estimator of  $\sigma^2$ . Therefore  $\hat{\boldsymbol{\Sigma}} = s^2(\mathbf{X}^\top \mathbf{X})^{-1}$  is the estimated covariance matrix as available in the regression printout. Therefore  $\hat{\mathbf{V}} = s^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top$  is the estimate of the covariance matrix of  $\mathbf{R}\hat{\boldsymbol{\beta}}$ . Another way to write (30.4.2) is therefore

$$(30.4.3) \quad B(\mathbf{y}) = \{\tilde{\mathbf{u}} \in \mathbb{R}^i : (\mathbf{R}\hat{\boldsymbol{\beta}} - \tilde{\mathbf{u}})^\top \hat{\mathbf{V}}^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \tilde{\mathbf{u}}) \leq iF_{(i, n-k; \alpha)}\}.$$

This formula allows a suggestive interpretation. whether  $\tilde{\mathbf{u}}$  lies in the confidence region or not depends on the Mahalanobis distance of the actual value of  $\mathbf{R}\hat{\boldsymbol{\beta}}$  would have from the distribution which  $\mathbf{R}\hat{\boldsymbol{\beta}}$  would have if the true parameter vector were to satisfy the constraint  $\mathbf{R}\boldsymbol{\beta} = \tilde{\mathbf{u}}$ . It is not the Mahalanobis distance itself but only an estimate of it because  $\sigma^2$  is replaced by its unbiased estimate  $s^2$ .

These formulas are also useful for drawing the confidence ellipses. The  $r$  which you need in equation (7.3.22) in order to draw the confidence ellipse is  $r = \sqrt{iF_{(i, n-k; \alpha)}}$ . This is the same as the local variable `mult` in the following S-function to draw this ellipse: its arguments are the center point (a 2-vector `d`), the estimated covariance matrix (a  $2 \times 2$  matrix `C`), the degrees of freedom in the denominator of the F-distribution (the scalar `df`), and the confidence level (the scalar `level` between 0 and 1 which defaults to 0.95 if not specified).

```
confelli <-
function(b, C, df, level = 0.95, xlab = "", ylab = "", add=T, prec=51)
```

```
# Plot an ellipse with "covariance matrix" C, center b, and P-content
# level according the F(2,df) distribution.
# Sent to S-NEWS on May 19, 1999 by Roger Koenker
# Department of Economics
# University of Illinois
# Champaign, IL 61820
# url: http://www.econ.uiuc.edu
# email roger@ysidro.econ.uiuc.edu
# vox: 217-333-4558
# fax: 217-244-6678.
# Included in the ecmnet package with his permission.
```

```
{
d <- sqrt(diag(C))
dfvec <- c(2, df)
phase <- acos(C[1, 2]/(d[1] * d[2]))
angles <- seq( - (PI), PI, len = prec)
mult <- sqrt(dfvec[1] * qf(level, dfvec[1], dfvec[2]))
xpts <- b[1] + d[1] * mult * cos(angles)
ypts <- b[2] + d[2] * mult * cos(angles + phase)
if(add) lines(xpts, ypts)
else plot(xpts, ypts, type = "l", xlab = xlab, ylab = ylab)
}
```

The mathematics why this works is in Problem 146.

PROBLEM 345. *3 points* In the regression model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  you observe  $\mathbf{y}$  and the (nonstochastic)  $\mathbf{X}$  and you construct the following confidence region  $B(\mathbf{y})$  for  $\mathbf{R}\boldsymbol{\beta}$ , where  $\mathbf{R}$  is a  $i \times k$  matrix with full row rank:

$$(30.4.4) \quad B(\mathbf{y}) = \{\mathbf{u} \in \mathbb{R}^i : (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}) \leq i s^2 F_{(i, n-k; \alpha)}\}.$$

Compute the probability that  $B$  contains the true  $\mathbf{R}\boldsymbol{\beta}$ .

ANSWER.

$$(30.4.5) \quad \Pr[B(\mathbf{y}) \ni \mathbf{R}\boldsymbol{\beta}] = \Pr[(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{R}\boldsymbol{\beta})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{R}\boldsymbol{\beta}) \leq iF_{(i, n-k; \alpha)} s^2] =$$

$$(30.4.6) \quad = \Pr\left[\frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{R}\boldsymbol{\beta})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{R}\boldsymbol{\beta}) / i}{s^2} \leq F_{(i, n-k; \alpha)}\right] = 1 - \alpha$$

This interpretation with the Mahalanobis distance is commonly used for the construction of  $t$ -Intervals. A  $t$ -interval is a special case of the above confidence region for the case  $i = 1$ . The confidence interval with confidence level  $1 - \alpha$  for the scalar parameter  $u = \mathbf{r}^\top \boldsymbol{\beta}$ , where  $\mathbf{r} \neq \mathbf{o}$  is a vector of constant coefficients, can be written as

$$(30.4.7) \quad B(\mathbf{y}) = \{u \in \mathbb{R} : |u - \mathbf{r}^\top \hat{\boldsymbol{\beta}}| \leq t_{(n-k; \alpha/2)} s_{\mathbf{r}^\top \hat{\boldsymbol{\beta}}}\}.$$

What do those symbols mean?  $\hat{\boldsymbol{\beta}}$  is the least squares estimator of  $\boldsymbol{\beta}$ .  $t_{(n-k; \alpha/2)}$  is the upper  $\alpha/2$ -quantile of the  $t$  distribution with  $n - k$  degrees of freedom, i.e. that scalar  $c$  for which a random variable  $t$  which has a  $t$  distribution with  $n - k$  degrees of freedom satisfies  $\Pr[t \geq c] = \alpha/2$ . Since by symmetry  $\Pr[t \leq -c] = \alpha/2$  as well, one obtains the inequality relevant for a two-sided test:

$$(30.4.8) \quad \Pr[|t| \geq t_{(n-k; \alpha/2)}] = \alpha.$$

Finally,  $s_{\mathbf{r}^\top \hat{\boldsymbol{\beta}}}$  is the estimated standard deviation of  $\mathbf{r}^\top \hat{\boldsymbol{\beta}}$ .



It is computed by the following three steps: First write down the variance of  $\mathbf{r}^\top \hat{\boldsymbol{\beta}}$ :

$$(30.4.9) \quad \text{var}[\mathbf{r}^\top \hat{\boldsymbol{\beta}}] = \sigma^2 \mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}.$$

Secondly, replace  $\sigma^2$  by its unbiased estimator  $s^2 = (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) / (n - k)$ , and thirdly take the square root. This gives  $s_{\mathbf{r}^\top \hat{\boldsymbol{\beta}}} = s \sqrt{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}}$ .

PROBLEM 346. Which element(s) on the right hand side of (30.4.7) depend(s) on  $\mathbf{y}$ ?

ANSWER.  $\hat{\boldsymbol{\beta}}$  depends on  $\mathbf{y}$ , and also  $s_{\mathbf{r}^\top \hat{\boldsymbol{\beta}}}$  depends on  $\mathbf{y}$  through  $s^2$ . □

Let us verify that the coverage probability, i.e., the probability that the confidence interval constructed using formula (30.4.7) contains the true value  $\mathbf{r}^\top \boldsymbol{\beta}$ , is, as claimed,  $1 - \alpha$ :

(30.4.10)

$$(30.4.11) \quad \begin{aligned} \Pr[\mathbf{B}(\mathbf{y}) \ni \mathbf{r}^\top \boldsymbol{\beta}] &= \Pr[|\mathbf{r}^\top \boldsymbol{\beta} - \mathbf{r}^\top \hat{\boldsymbol{\beta}}| \leq t_{(n-k; \alpha/2)} s_{\mathbf{r}^\top \hat{\boldsymbol{\beta}}}] \\ &= \Pr\left[|\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}| \leq t_{(n-k; \alpha/2)} s \sqrt{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}}\right] \end{aligned}$$

$$(30.4.12) \quad = \Pr\left[\left|\frac{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}}{s \sqrt{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}}}\right| \leq t_{(n-k; \alpha/2)}\right]$$

$$(30.4.13) \quad = \Pr\left[\left|\frac{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}}{\sigma \sqrt{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}}}\right| \Big/ \frac{s}{\sigma} \leq t_{(n-k; \alpha/2)}\right] = 1 - \alpha,$$

This last equality holds because the expression left of the big slash is a standard normal, and the expression on the right of the big slash is the square root of an independent  $\chi_{n-k}^2$  divided by  $n - k$ . The random variable between the absolute signs has therefore a t-distribution, and (30.4.13) follows from (30.4.8).

In R, one obtains  $t_{(n-k; \alpha/2)}$  by giving the command `qt(1-alpha/2, n-p)`. Here `qt` stands for t-quantile [BCW96, p. 48]. One needs `1-alpha/2` instead of `alpha/2` because it is the usual convention for quantiles (or cumulative distribution functions) to be defined as lower quantiles, i.e., as the probabilities of a random variable being  $\leq$  a given number, while test statistics are usually designed in such a way that the significant values are the high values, i.e., for testing one needs the upper quantiles.

There is a basic duality between confidence intervals and hypothesis tests. Chapter 31 is therefore a discussion of the same subject under a slightly different angle:

## CHAPTER 31

## Three Principles for Testing a Linear Constraint

We work in the model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with normally distributed errors  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ . There are three basic approaches to test the null hypothesis  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$ . In the linear model, these three approaches are mathematically equivalent, but if one goes over to nonlinear least squares or maximum likelihood estimators, they lead to different (although asymptotically equivalent) tests.

(1) (“Wald Criterion”) Compute the vector of OLS estimates  $\hat{\boldsymbol{\beta}}$ , and reject the null hypothesis if  $\mathbf{R}\hat{\boldsymbol{\beta}}$  is “too far away” from  $\mathbf{u}$ . For this criterion one only needs the unconstrained estimator, not the constrained one.

(2) (“Likelihood Ratio Criterion”) Estimate the model twice: once with the constraint  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$ , and once without the constraint. Reject the null hypothesis if the model with the constraint imposed has a much worse fit than the model without the constraint.

(3) (“Lagrange Multiplier Criterion”) This third criterion is based on the constrained estimator only. It has two variants. In its “score test” variant, one rejects the null hypothesis if the vector of derivatives of the unconstrained least squares objective function, evaluated at the constrained estimate  $\hat{\boldsymbol{\beta}}$ , is too far away from  $\mathbf{0}$ . In the variant which has given this Criterion its name, one rejects if the vector of Lagrange multipliers needed for imposing the constraint is too far away from  $\mathbf{0}$ .

Many textbooks inadvertently and implicitly distinguish between (1) and (2) as follows: they introduce the t-test for one parameter by principle (1), and the F-test for several parameters by principle (2). Later, the student is surprised to find out that the t-test and the F-test in one dimension are equivalent, i.e., that the difference between t-test and F-test has nothing to do with the dimension of the parameter vector to be tested. Some textbooks make the distinction between (1) and (2) *explicit*. For instance [Chr87, p. 29ff] distinguishes between “testing linear parametric functions” and “testing models.” However the distinction between all 3 principles has been introduced into the linear model only after the discovery that these three principles give different but asymptotically equivalent tests in the Maximum Likelihood estimation. Compare [DM93, Chapter 3.6] about this.

## 31.1. Mathematical Detail of the Three Approaches

(1) For the “Wald criterion” we must specify what it means that  $\mathbf{R}\hat{\boldsymbol{\beta}}$  is “too far away” from  $\mathbf{u}$ . The Mahalanobis distance gives such a criterion: If the true  $\boldsymbol{\beta}$  satisfies  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$ , then  $\mathbf{R}\hat{\boldsymbol{\beta}} \sim (\mathbf{u}, \sigma^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)$ , and the Mahalanobis distance of the observed value of  $\mathbf{R}\hat{\boldsymbol{\beta}}$  from this distribution is a logical candidate for the Wald criterion. The only problem is that  $\sigma^2$  is not known, therefore we have to use the “studentized” Mahalanobis distance in which  $\sigma^2$  is replaced by  $s^2$ . Conventionally, in the context of linear regression, the Mahalanobis distance is also divided by the number of degrees of freedom; this normalizes its expected value to 1. Replacing  $\sigma^2$  by  $s^2$  and dividing by  $i$  gives the test statistic

$$(31.1.1) \quad \frac{1}{i} \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})}{s^2}.$$

(2) Here are the details for the second approach, the “goodness-of-fit criterion”. In order to compare the fit of the models, we look at the attained SSE’s. Of course the constrained SSE<sub>r</sub> is always larger than the unconstrained SSE<sub>u</sub>, even if the true parameter vector satisfies the constraint. But if we divide SSE<sub>r</sub> by its degrees of freedom  $n + i - k$ , it is an unbiased estimator of  $\sigma^2$  if the constraint holds and it is biased upwards if the constraint does not hold. The unconstrained SSE<sub>u</sub>, divided by its degrees of freedom, on the other hand, is always an unbiased estimator of  $\sigma^2$ . If the constraint holds, the SSE’s divided by their respective degrees of freedom should give roughly equal numbers. According to this, a feasible test statistic would be

$$(31.1.2) \quad \frac{\text{SSE}_r / (n + i - k)}{\text{SSE}_u / (n - k)}$$

and one would reject if this is too much  $> 1$ . The following variation of this is more convenient, since its distribution does not depend on  $n$ ,  $k$  and  $i$  separately, but only through  $n - k$  and  $i$ .

$$(31.1.3) \quad \frac{(\text{SSE}_r - \text{SSE}_u) / i}{\text{SSE}_u / (n - k)}$$

It still has the property that the numerator is an unbiased estimator of  $\sigma^2$  if the constraint holds and biased upwards if the constraint does not hold, and the denominator is always an unbiased estimator. Furthermore, in this variation, the numerator and denominator are independent random variables. If this test statistic is much larger than 1, then the constraints are incompatible with the data and the null hypothesis must be rejected. The statistic (31.1.3) can also be written as

$$(31.1.4) \quad \frac{(\text{SSE}_{\text{constrained}} - \text{SSE}_{\text{unconstrained}}) / \text{number of constraints}}{\text{SSE}_{\text{unconstrained}} / (\text{numb. of observations} - \text{numb. of coefficients in unconstr. model})}$$

The equivalence of formulas (31.1.1) and (31.1.4) is a simple consequence of (22.7.2).

(3) And here are the details about the score test variant of the Lagrange multiplier criterion: The Jacobian of the least squares objective function is

$$(31.1.5) \quad \frac{\partial}{\partial \boldsymbol{\beta}^\top} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = -2(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{X}.$$

This is a row vector consisting of all the partial derivatives. Taking its transpose, in order to get a column vector, and plugging the constrained least squares estimate  $\hat{\boldsymbol{\beta}}$  into it gives  $-2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ . Again we need the Mahalanobis distance of this observed value from the distribution which the random variable

$$(31.1.6) \quad -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

has if the true  $\boldsymbol{\beta}$  satisfies  $\mathbf{R}\boldsymbol{\beta} = \mathbf{u}$ . If this constraint is satisfied,  $\hat{\boldsymbol{\beta}}$  is unbiased, therefore (31.1.6) has expected value zero. Furthermore, if one premultiplies (22.7.1) by  $\mathbf{X}^\top$  one gets  $\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})$ , therefore  $\nu[\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})] = \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \mathbf{R}$ ; and now one can see that  $\frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{X})^{-1}$  is a g-inverse of this covariance matrix. Therefore the Mahalanobis distance of the observed value from the distribution is

$$(31.1.7) \quad \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

The Lagrange multiplier statistic is based on the restricted estimator alone. If one wanted to take this principle seriously one would have to replace  $\sigma^2$  by the unbiased estimate from the restricted model to get the “score form” of the Lagrange Multiplier Test statistic. But in the linear model this leads to it that the denominator in the test statistic is no longer independent of the numerator, and since the test statistic as a function of the ratio of the constrained and unconstrained estimates of  $\sigma^2$  anyway, one will only get yet another monotonic transformation of the same test statistic. If one were to use the unbiased estimate from the unrestricted model, one would exactly get the Wald statistic back, as one can verify using (22.3.13).

This same statistic can also be motivated in terms of the Lagrange multipliers, and this is where this testing principle has its name from, although the applications usually use the score form. According to (22.3.12), the Lagrange multiplier is  $\boldsymbol{\lambda} = 2(\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})$ . If the constraint holds, then  $\boldsymbol{\varepsilon}[\boldsymbol{\lambda}] = \mathbf{o}$ , and  $\nu[\boldsymbol{\lambda}] = 4\sigma^2 (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1}$ . The Mahalanobis distance of the observed value from this distribution is

$$(31.1.8) \quad \boldsymbol{\lambda}^\top (\nu[\boldsymbol{\lambda}])^{-1} \boldsymbol{\lambda} = \frac{1}{4\sigma^2} \boldsymbol{\lambda}^\top \mathbf{R} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}$$

Using (22.7.1) one can verify that this is the same as (31.1.7).

PROBLEM 347. Show that (31.1.7) is equal to the righthand side of (31.1.8).

PROBLEM 348. 10 points Prove that  $\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} - \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}$  can be written alternatively the following five ways:

$$(31.1.9) \quad \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} - \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})$$

$$(31.1.10) \quad = (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u})$$

$$(31.1.11) \quad = \frac{1}{4} \boldsymbol{\lambda}^\top \mathbf{R} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda}$$

$$(31.1.12) \quad = \hat{\boldsymbol{\varepsilon}}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\boldsymbol{\varepsilon}}$$

$$(31.1.13) \quad = (\hat{\boldsymbol{\varepsilon}} - \hat{\boldsymbol{\varepsilon}})^\top (\hat{\boldsymbol{\varepsilon}} - \hat{\boldsymbol{\varepsilon}})$$

Furthermore show that

$$(31.1.14) \quad \mathbf{X}^\top \mathbf{X} \text{ is } \sigma^2 \text{ times a g-inverse of } \nu[\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}]$$

$$(31.1.15) \quad (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \text{ is } \sigma^2 \text{ times the inverse of } \nu[\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}]$$

$$(31.1.16) \quad \frac{1}{4} \mathbf{R} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \text{ is } \sigma^2 \text{ times the inverse of } \nu[\boldsymbol{\lambda}]$$

$$(31.1.17) \quad (\mathbf{X}^\top \mathbf{X})^{-1} \text{ is } \sigma^2 \text{ times a g-inverse of } \nu[\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})]$$

$$(31.1.18) \quad \mathbf{I} \text{ is } \sigma^2 \text{ times a g-inverse of } \nu[\hat{\boldsymbol{\varepsilon}} - \hat{\boldsymbol{\varepsilon}}]$$

and show that  $-2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$  is the gradient of the SSE objective function evaluated at  $\hat{\boldsymbol{\beta}}$ . By the way, one should be a little careful in interpreting (31.1.12) because  $\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is not  $\sigma^2$  times the g-inverse of  $\nu[\hat{\boldsymbol{\varepsilon}}]$ .

ANSWER.

$$(31.1.19) \quad \hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + \hat{\boldsymbol{\varepsilon}},$$

and since  $\mathbf{X}^\top \hat{\boldsymbol{\varepsilon}} = \mathbf{o}$ , the righthand decomposition is an orthogonal decomposition. This gives (31.1.9) above:

$$(31.1.20) \quad \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}},$$

Using (22.3.13) one obtains  $\nu[\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} \mathbf{R} (\mathbf{X}^\top \mathbf{X})^{-1}$ . This is a singular matrix, and one verifies immediately that  $\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$  is a g-inverse of it.

To obtain (31.1.10), which is (22.7.2), one has to plug (22.3.13) into (31.1.20). Clearly,  $\nu[\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{u}] = \sigma^2 \mathbf{R} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top$ .

For (31.1.11) one needs the formula for the Lagrange multiplier (22.3.12).

The test statistic defined alternatively either by (31.1.1) or (31.1.4) or (31.1.8) or (31.1.8) has the following nice properties:

- $E(\text{SSE}_u) = E(\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}) = \sigma^2(n - k)$ , which holds whether or not the constraint is true. Furthermore it was shown earlier that

$$(31.1.21) \quad E(\text{SSE}_r - \text{SSE}_u) = \sigma^2 i + (\mathbf{R}\boldsymbol{\beta} - \mathbf{u})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\boldsymbol{\beta} - \mathbf{u}),$$

i.e., this expected value is equal to  $\sigma^2 i$  if the constraint is true, and larger otherwise. If one divides  $\text{SSE}_u$  and  $\text{SSE}_r - \text{SSE}_u$  by their respective degrees of freedom, as is done in (31.1.4), one obtains therefore: the denominator is always an unbiased estimator of  $\sigma^2$ , regardless of whether the null hypothesis is true or not. The numerator is an unbiased estimator of  $\sigma^2$  when the null hypothesis is correct, and has a *positive* bias otherwise.

- If the distribution of  $\boldsymbol{\varepsilon}$  is normal, then numerator and denominator are independent. The numerator is a function of  $\hat{\boldsymbol{\beta}}$  and the denominator one of  $\hat{\boldsymbol{\varepsilon}}$ , and  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\varepsilon}}$  are independent.
- Again under assumption of normality, numerator and denominator are distributed as  $\sigma^2 \chi^2$  with  $i$  and  $n - k$  degrees of freedom, divided by their respective degrees of freedom. If one divides them, the common factor  $\sigma^2$  cancels out, and the ratio has a F distribution. Since both numerator and denominator have the same expected value  $\sigma^2$ , the value of this F distribution should be in the order of magnitude of 1. If it is much larger than that, the null hypothesis is to be rejected. (Precise values in the F-tables).

### 31.2. Examples of Tests of Linear Hypotheses

Some tests can be read off directly from the computer printouts. One example is the t-tests for an individual component of  $\boldsymbol{\beta}$ . The situation is  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where  $\boldsymbol{\beta} = [\beta_1 \ \cdots \ \beta_k]^\top$ , and we want to test  $\beta_j = u$ . Here  $\mathbf{R} = \mathbf{e}_j = [0 \ \cdots \ 0 \ 1 \ 0 \ \cdots \ 0]$ , with the 1 on the  $j$ th place, and  $\mathbf{u}$  is the 1-vector  $u$ , and  $i = 1$ . Therefore  $\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top = d_{jj}$ , the  $j$ th diagonal element of  $(\mathbf{X}^\top \mathbf{X})^{-1}$ , and (31.1.1) becomes

$$(31.2.1) \quad \frac{(\hat{\beta}_j - u)^2}{s^2 d_{jj}} \sim F_{1, n-k} \quad \text{when } H \text{ is true.}$$

This is the square of a random variable which has a t-distribution:

$$(31.2.2) \quad \frac{\hat{\beta}_j - u}{s \sqrt{d_{jj}}} \sim t_{n-k} \quad \text{when } H \text{ is true.}$$

This latter test statistic is simply  $\hat{\beta}_j - u$  divided by the estimated standard deviation of  $\hat{\beta}_j$ .

If one wants to test that a certain linear combination of the parameter values is equal to (or bigger than or smaller than) a given value, say  $\mathbf{r}^\top \boldsymbol{\beta} = u$ , one can use a

t-test as well. The test statistic is, again, simply  $\mathbf{r}^\top \hat{\boldsymbol{\beta}} - u$  divided by the estimated standard deviation of  $\mathbf{r}^\top \hat{\boldsymbol{\beta}}$ :

$$(31.2.3) \quad \frac{\mathbf{r}^\top \hat{\boldsymbol{\beta}} - u}{s \sqrt{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}}} \sim t_{n-k} \quad \text{when } H \text{ is true.}$$

By this one can for instance also test whether the sum of certain regression coefficients is equal to 1, or whether two regression coefficients are equal to each other (but not the hypothesis that three coefficients are equal to each other).

Many textbooks use the Wald criterion to derive the t-test, and the Likelihood Ratio criterion to derive the F-test. Our approach showed that the Wald criterion can be used for simultaneous testing of several hypotheses as well. The t-test is equivalent to an F-test if only one hypothesis is tested, i.e., if  $\mathbf{R}$  is a row vector. The only difference is that with the t-test one can test one-sided hypotheses, with the F-test one cannot.

Next let us discuss the test for the existence of a relationship, “the” F-test which every statistics package performs automatically whenever the regression has a constant term: it is the test whether all the slope parameters are zero, such that only the intercept may take a nonzero value.

**PROBLEM 349.** 4 points In the model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with intercept, show that the test statistic for testing whether all the slope parameters are zero is

$$(31.2.4) \quad \frac{(\mathbf{y}^\top \mathbf{X} \hat{\boldsymbol{\beta}} - n\bar{y}^2)/(k - 1)}{(\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \hat{\boldsymbol{\beta}})/(n - k)}$$

This is [Seb77, equation (4.26) on p. 110]. What is the distribution of this test statistic if the null hypothesis is true (i.e., if all the slope parameters are zero)?

**ANSWER.** The distribution is  $\sim F_{k-1, n-k}$ . (31.2.4) is most conveniently derived from (31.1.1). In the constrained model, which has only a constant term and no other explanatory variables,  $\mathbf{y} = \boldsymbol{\nu}\mu + \boldsymbol{\varepsilon}$ , the BLUE is  $\hat{\mu} = \bar{y}$ . Therefore the constrained residual sum of squares  $\text{SSE}_{\text{const.}}$  is what is commonly called SST (“total” or, more precisely, “corrected total” sum of squares):

$$(31.2.5) \quad \text{SSE}_{\text{const.}} = \text{SST} = (\mathbf{y} - \boldsymbol{\nu}\bar{y})^\top (\mathbf{y} - \boldsymbol{\nu}\bar{y}) = \mathbf{y}^\top (\mathbf{y} - \boldsymbol{\nu}\bar{y}) = \mathbf{y}^\top \mathbf{y} - n\bar{y}^2$$

while the unconstrained residual sum of squares is what is usually called SSE:

$$(31.2.6) \quad \text{SSE}_{\text{unconst.}} = \text{SSE} = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{y}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \hat{\boldsymbol{\beta}}.$$

This last equation because  $\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{X}^\top \hat{\boldsymbol{\varepsilon}} = \mathbf{o}$ . A more elegant way is perhaps

$$(31.2.7) \quad \text{SSE}_{\text{unconst.}} = \text{SSE} = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = \mathbf{y}^\top \mathbf{M}^\top \mathbf{M} \mathbf{y} = \mathbf{y}^\top \mathbf{M} \mathbf{y} = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \hat{\boldsymbol{\beta}}$$

According to (14.3.12) we can write  $SSR = SST - SSE$ , therefore the F-statistic is

$$(31.2.8) \quad \frac{SSR/(k-1)}{SSE/(n-k)} = \frac{(\mathbf{y}^\top \mathbf{X} \hat{\boldsymbol{\beta}} - n\bar{y}^2)/(k-1)}{(\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \hat{\boldsymbol{\beta}})/(n-k)} \sim F_{k-1, n-k} \quad \text{if } H_0 \text{ is true.}$$

□

**PROBLEM 350.** *2 points* Can one compute the value of the F-statistic testing for the existence of a relationship if one only knows the coefficient of determination  $R^2 = SSR/SST$ , the number of observations  $n$ , and the number of regressors (counting the constant term as one of the regressors)  $k$ ?

ANSWER.

$$(31.2.9) \quad F = \frac{SSR/(k-1)}{SSE/(n-k)} = \frac{n-k}{k-1} \frac{SSR}{SST - SSR} = \frac{n-k}{k-1} \frac{R^2}{1 - R^2}.$$

□

Other, similar F-tests are: the F-test that all among a number of additional variables have the coefficient zero, the F-test that three or more coefficients are equal. One can use the t-test for testing whether two coefficients are equal, but not for three. It may be possible that the t-test for  $\beta_1 = \beta_2$  does not reject and the t-test for  $\beta_2 = \beta_3$  does not reject either, but the t-test for  $\beta_1 = \beta_3$  does reject!

**PROBLEM 351.** *4 points* [Seb77, exercise 4b.5 on p. 109/10] In the model  $\mathbf{y} = \boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon} \sim N(\mathbf{o}, \sigma^2 \mathbf{I})$  and subject to the constraint  $\boldsymbol{\iota}^\top \boldsymbol{\beta} = 0$ , which we had in Problem 291, compute the test statistic for the hypothesis  $\beta_1 = \beta_3$ .

ANSWER. In this problem, the “unconstrained” model for the purposes of testing is already constrained, it is subject to the constraint  $\boldsymbol{\iota}^\top \boldsymbol{\beta} = 0$ . The “constrained” model has the additional constraint  $\mathbf{R}\boldsymbol{\beta} = \begin{bmatrix} 1 & 0 & -1 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} = 0$ . In Problem 291 we computed the “unconstrained” estimates  $\hat{\boldsymbol{\beta}} = \mathbf{y} - \boldsymbol{\iota}\bar{y}$  and  $s^2 = n\bar{y}^2 = (\mathbf{y}_1 + \cdots + \mathbf{y}_n)^2/n$ . You are allowed to use this without proving it again. Therefore  $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{y}_1 - \mathbf{y}_3$ ; its variance is  $2\sigma^2$ , and the F test statistic is  $\frac{n(\mathbf{y}_1 - \mathbf{y}_3)^2}{2(\mathbf{y}_1 + \cdots + \mathbf{y}_n)^2} \sim F_{1,1}$ . The “unconstrained” model had 4 parameters subject to one constraint, therefore it had 3 free parameters, i.e.,  $k = 3$ ,  $n = 4$ , and  $j = 1$ . □

Another important F-test is the “Chow test” named by its popularizer Chow [Cho60]: it tests whether two regressions have equal coefficients (assuming that the disturbance variances are equal). For this one has to run three regressions. If the first regression has  $n_1$  observations and sum of squared error  $SSE_1$ , and the second regression  $n_2$  observations and  $SSE_2$ , and the combined regression (i.e., the restricted model) has  $SSE_r$ , then the test statistic is

$$(31.2.10) \quad \frac{(SSE_r - SSE_1 - SSE_2)/k}{(SSE_1 + SSE_2)/(n_1 + n_2 - 2k)}.$$

If  $n_2 < k$ , the second regression cannot be run by itself. In this case, the unconstrained model has “too many” parameters: they give an exact fit for the second group of observations  $SSE_2 = 0$ , and in addition not all parameters are identifiable. In effect this second regression has only  $n_2$  parameters. These parameters can be considered dummy variables for every observation, i.e., this test can be interpreted to be a test whether the  $n_2$  additional observations come from the same population as the  $n_1$  first ones. The test statistic becomes

$$(31.2.11) \quad \frac{(SSE_r - SSE_1)/n_2}{SSE_1/(n_1 - k)}.$$

This latter is called the “predictive Chow test,” because in its Wald version it looks at the prediction errors involving observations in the second regression.

The following is a special case of the Chow test, in which one can give a simple formula for the test statistic.

**PROBLEM 352.** Assume you have  $n_1$  observations  $\mathbf{u}_j \sim N(\mu_1, \sigma^2)$  and  $n_2$  observations  $\mathbf{v}_j \sim N(\mu_2, \sigma^2)$ , all independent of each other, and you want to test whether  $\mu_1 = \mu_2$ . (Note that the variances are known to be equal).

- a. *2 points* Write the model in the form  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ .

ANSWER.

$$(31.2.12) \quad \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\iota}_1 \mu_1 + \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\iota}_2 \mu_2 + \boldsymbol{\varepsilon}_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\iota}_1 & \mathbf{o} \\ \mathbf{o} & \boldsymbol{\iota}_2 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\varepsilon}_2 \end{bmatrix}.$$

here  $\boldsymbol{\iota}_1$  and  $\boldsymbol{\iota}_2$  are vectors of ones of appropriate lengths.

- b. *2 points* Compute  $(\mathbf{X}^\top \mathbf{X})^{-1}$  in this case.

ANSWER.

$$(31.2.13) \quad \mathbf{X}^\top \mathbf{X} = \begin{bmatrix} \boldsymbol{\iota}_1^\top & \mathbf{o}^\top \\ \mathbf{o}^\top & \boldsymbol{\iota}_2^\top \end{bmatrix} \begin{bmatrix} \boldsymbol{\iota}_1 & \mathbf{o} \\ \mathbf{o} & \boldsymbol{\iota}_2 \end{bmatrix} = \begin{bmatrix} n_1 & 0 \\ 0 & n_2 \end{bmatrix}$$

$$(31.2.14) \quad (\mathbf{X}^\top \mathbf{X})^{-1} = \begin{bmatrix} \frac{1}{n_1} & 0 \\ 0 & \frac{1}{n_2} \end{bmatrix}$$

- c. *2 points* Compute  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  in this case.

ANSWER.

$$(31.2.15) \quad \mathbf{X}^\top \mathbf{y} = \begin{bmatrix} \boldsymbol{\iota}_1^\top & \mathbf{o}^\top \\ \mathbf{o}^\top & \boldsymbol{\iota}_2^\top \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{n_1} u_i \\ \sum_{j=1}^{n_2} v_j \end{bmatrix}$$

$$(31.2.16) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \begin{bmatrix} \frac{1}{n_1} & 0 \\ 0 & \frac{1}{n_2} \end{bmatrix} \begin{bmatrix} \sum_{i=1}^{n_1} u_i \\ \sum_{j=1}^{n_2} v_j \end{bmatrix} = \begin{bmatrix} \bar{u} \\ \bar{v} \end{bmatrix}$$

• d. 3 points Compute  $SSE = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$  and  $s^2$ , the unbiased estimator of  $\sigma^2$ , in this case.

ANSWER.

$$(31.2.17) \quad \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} - \begin{bmatrix} \iota_1 & \mathbf{o} \\ \mathbf{o} & \iota_2 \end{bmatrix} \begin{bmatrix} \bar{\mathbf{u}} \\ \bar{\mathbf{v}} \end{bmatrix} = \begin{bmatrix} \mathbf{u} - \iota_1 \bar{\mathbf{u}} \\ \mathbf{v} - \iota_2 \bar{\mathbf{v}} \end{bmatrix}$$

$$(31.2.18) \quad SSE = \sum_{i=1}^{n_1} (u_i - \bar{u})^2 + \sum_{j=1}^{n_2} (v_j - \bar{v})^2$$

$$(31.2.19) \quad s^2 = \frac{\sum_{i=1}^{n_1} (u_i - \bar{u})^2 + \sum_{j=1}^{n_2} (v_j - \bar{v})^2}{n_1 + n_2 - 2}$$

□

• e. 1 point Next, the hypothesis  $\mu_1 = \mu_2$  must be written in the form  $\mathbf{R}\boldsymbol{\beta} = u$ . Since in the present case  $\mathbf{R}$  has just has one row, it should be written as a row-vector  $\mathbf{R} = \mathbf{r}^\top$ , and since the vector  $\mathbf{u}$  has only one component, it should be written as a scalar  $u$ , i.e., the hypothesis should be written in the form  $\mathbf{r}^\top \boldsymbol{\beta} = u$ . What are  $\mathbf{r}$  and  $u$  in our case?

ANSWER. Since  $\boldsymbol{\beta} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ , the constraint can be written as

$$(31.2.20) \quad \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = 0 \quad \text{i.e.,} \quad \mathbf{r} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad \text{and} \quad u = 0$$

□

• f. 2 points Compute the standard deviation of  $\mathbf{r}^\top \hat{\boldsymbol{\beta}}$ .

ANSWER. First compute the variance and then take the square root.

$$(31.2.21) \quad \text{var}[\mathbf{r}^\top \hat{\boldsymbol{\beta}}] = \sigma^2 \mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r} = \sigma^2 \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} \frac{1}{n_1} & 0 \\ 0 & \frac{1}{n_2} \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \sigma^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)$$

One can also see this without matrix algebra.  $\text{var}[\bar{u} = \sigma^2 \frac{1}{n_1}]$ ,  $\text{var}[\bar{v} = \sigma^2 \frac{1}{n_2}]$ , and since  $\bar{u}$  and  $\bar{v}$  are independent, the variance of the difference is the sum of the variances. □

• g. 2 points Use (31.2.3) to derive the formula for the t-test.

ANSWER. The test statistic is  $\bar{u} - \bar{v}$  divided by its estimated standard deviation, i.e.,

$$(31.2.22) \quad \frac{\bar{u} - \bar{v}}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2} \quad \text{when } H \text{ is true.}$$

□

PROBLEM 353. [Seb77, exercise 4d-3] Given  $n + 1$  observations  $y_j$  from a  $N(\mu, \sigma^2)$ . After the first  $n$  observations, it is suspected that a sudden change in the mean of the distribution occurred, i.e., that  $\mathbf{y}_{n+1} \sim N(\nu, \sigma^2)$  with  $\nu \neq \mu$ . We will

use here three different approaches to derive the same test statistic for testing hypothesis that the  $n + 1$ st observation has the same population mean as the previous observations, i.e., that  $\nu = \mu$ , against the two-sided alternative. The formulas for this statistic should be given in terms of the observations  $y_i$ . It is recommended use the notation  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  and  $\bar{y} = \frac{1}{n+1} \sum_{j=1}^{n+1} y_j$ .

• a. 3 points First you should derive this statistic by testing whether  $\nu - \mu = 0$  (the “Wald principle”). For this you must compute the BLUE of  $\nu - \mu$  and its standard deviation and construct the t statistic from this.

ANSWER. BLUE of  $\mu$  is  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ , and that of  $\nu$  is  $y_{n+1}$ . BLUE of  $\nu - \mu$  is  $\bar{y} - y_{n+1}$ . Because of independence  $\text{var}[\bar{y} - y_{n+1}] = \text{var}[\bar{y}] + \text{var}[y_{n+1}] = \sigma^2((1/n) + 1) = \sigma^2(n + 1)/n$ . Standard deviation is  $\sigma \sqrt{(n + 1)/n}$ .

For the denominator in the t-statistic you need the  $s^2$  from the unconstrained regression, which is

$$(31.2.23) \quad s^2 = \frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y})^2$$

What happened to the  $(n + 1)$ st observation here? It always has a zero residual. And the factor  $1/(n - 1)$  should really be written  $1/(n + 1 - 2)$ : there are  $n + 1$  observations and 2 parameters.

Divide  $\bar{y} - y_{n+1}$  by its standard deviation and replace  $\sigma$  by  $s$  (the square root of  $s^2$ ) to get the t statistic

$$(31.2.24) \quad \frac{\bar{y} - y_{n+1}}{s \sqrt{1 + \frac{1}{n}}}$$

• b. 2 points One can interpret this same formula also differently (and this is why this test is sometimes called the “predictive” Chow test). Compute the BLUE of  $\mathbf{y}_{n+1}$  on the basis of the first  $n$  observations, call it  $\hat{\mathbf{y}}(n + 1)_{n+1}$ . Show that the predictive residual  $y_{n+1} - \hat{\mathbf{y}}(n + 1)_{n+1}$ , divided by the square root of  $\text{MSE}[\hat{\mathbf{y}}(n + 1)_{n+1}; \mathbf{y}_{n+1}]$ , with  $\sigma$  replaced by  $s$  (based on the first  $n$  observations only), is equal to the above t statistic.

ANSWER. BLUP of  $\mathbf{y}_{n+1}$  based on first  $n$  observations is  $\bar{y}$  again. Since it is unbiased  $\text{MSE}[\bar{y}; \mathbf{y}_{n+1}] = \text{var}[\bar{y} - y_{n+1}] = \sigma^2(n + 1)/n$ . From now on everything is as in part a.

• c. 6 points Next you should show that the above two formulas are identical to the statistic based on comparing the SSEs of the constrained and unconstrained models (the likelihood ratio principle). Give a formula for the constrained SSE<sub>r</sub>, the unconstrained SSE<sub>u</sub>, and the F-statistic.

ANSWER. According to the Likelihood Ratio principle, one has to compare the residual sum of squares in the regressions under the assumption that the mean did not change with that under the assumption that the mean changed. If the mean did not change (constrained model), then  $\bar{y}$  is

OLS of  $\mu$ . In order to make it easier to derive the difference between constrained and unconstrained SSE, we will write the constrained SSE as follows:

$$SSE_r = \sum_{j=1}^{n+1} (y_j - \bar{y})^2 = \sum_{j=1}^{n+1} y_j^2 - (n+1)\bar{y}^2 = \sum_{j=1}^{n+1} y_j^2 - \frac{1}{n+1} (n\bar{y} + y_{n+1})^2$$

If one allows the mean to change (unconstrained model), then  $\bar{y}$  is the BLUE of  $\mu$ , and  $y_{n+1}$  is the BLUE of  $\nu$ .

$$SSE_u = \sum_{j=1}^n (y_j - \bar{y})^2 + (y_{n+1} - y_{n+1})^2 = \sum_{j=1}^n y_j^2 - n\bar{y}^2.$$

Now subtract:

$$\begin{aligned} SSE_r - SSE_u &= y_{n+1}^2 + n\bar{y}^2 - \frac{1}{n+1} (n\bar{y} + y_{n+1})^2 \\ &= y_{n+1}^2 + n\bar{y}^2 - \frac{1}{n+1} (n^2\bar{y}^2 + 2n\bar{y}y_{n+1} + y_{n+1}^2) \\ &= (1 - \frac{1}{n+1})y_{n+1}^2 + (n - \frac{n^2}{n+1})\bar{y}^2 - \frac{n}{n+1} 2\bar{y}y_{n+1} \\ &= \frac{n}{n+1} (y_{n+1} - \bar{y})^2. \end{aligned}$$

Interestingly, this depends on the first  $n$  observations only through  $\bar{y}$ .

Since the unconstrained model has  $n+1$  observations and 2 parameters, the test statistic is

$$(31.2.25) \quad \frac{SSE_r - SSE_u}{SSE_u / (n+1-2)} = \frac{\frac{n}{n+1} (y_{n+1} - \bar{y})^2}{\sum_1^n (y_j - \bar{y})^2 / (n-1)} = \frac{(y_{n+1} - \bar{y})^2 n(n-1)}{\sum_1^n (y_j - \bar{y})^2 (n+1)} \sim F_{1, n-1}$$

This is the square of the t statistic (31.2.24). □

### 31.2.1. Goodness of Fit Test.

PROBLEM 354. [Seb77, pp. 117–119] *Given a regression model with  $k$  independent variables. There are  $n$  observations of the vector of independent variables, and for each of these  $n$  values there is not one but  $r > 1$  different replicated observations of the dependent variable. This model can be written*

$$(31.2.26) \quad \mathbf{y}_{mq} = \sum_{j=1}^k x_{mj} \beta_j + \varepsilon_{mq} \quad \text{or} \quad \mathbf{y}_{mq} = \mathbf{x}_m^\top \boldsymbol{\beta} + \varepsilon_{mq},$$

where  $m = 1, \dots, n$ ,  $j = 1, \dots, k$ ,  $q = 1, \dots, r$ , and  $\mathbf{x}_m^\top$  is the  $m$ th row of the  $\mathbf{X}$ -matrix. For simplicity we assume that  $r$  does not depend on  $m$ , each observation of the independent variables has the same number of repetitions. We also assume that the  $n \times k$  matrix  $\mathbf{X}$  has full column rank.

• a. 2 points In this model it is possible to test whether the regression line is in fact a straight line. If it is not a straight line, then each observation of the dependent

variables  $\mathbf{x}_m$  has a different coefficient vector  $\boldsymbol{\beta}_m$  associated with it, i.e., the model is

$$(31.2.27) \quad \mathbf{y}_{mq} = \sum_{j=1}^k x_{mj} \beta_{mj} + \varepsilon_{mq} \quad \text{or} \quad \mathbf{y}_{mq} = \mathbf{x}_m^\top \boldsymbol{\beta}_m + \varepsilon_{mq}.$$

This unconstrained model does not have enough information to estimate any of individual coefficients  $\beta_{mj}$ . Explain how it is nevertheless still possible to compute  $SSE_u$ .

ANSWER. Even though the individual coefficients  $\beta_{mj}$  are not identified, their linear combination  $\eta_m = \mathbf{x}_m^\top \boldsymbol{\beta}_m = \sum_{j=1}^k x_{mj} \beta_{mj}$  is identified; one unbiased estimator, although by far not the best one, is any individual observation  $y_{mq}$ . This linear combination is all one needs to compute  $SSE_u$ , the sum of squared errors in the unconstrained model.

• b. 2 points Writing your estimate of  $\eta_m = \mathbf{x}_m^\top \boldsymbol{\beta}_m$  as  $\tilde{\eta}_m$ , give the formula of the sum of squared errors of this estimate, and by taking the first order conditions, show that the unconstrained least squares estimate of  $\eta_m$  is  $\hat{\eta}_m = \bar{y}_m$  for  $m = 1, \dots, n$  where  $\bar{y}_m = \frac{1}{r} \sum_{q=1}^r y_{mq}$  (i.e., the dot in the subscript indicates taking the mean).

ANSWER. If we know the  $\tilde{\eta}_m$  the sum of squared errors no longer depends on the independent observations  $\mathbf{x}_m$  but is simply

$$(31.2.28) \quad SSE_u = \sum_{m,q} (y_{mq} - \tilde{\eta}_m)^2$$

First order conditions are

$$(31.2.29) \quad \frac{\partial}{\partial \tilde{\eta}_h} \sum_{m,q} (y_{mq} - \tilde{\eta}_m)^2 = \frac{\partial}{\partial \tilde{\eta}_h} \sum_q (y_{hq} - \tilde{\eta}_h)^2 = -2 \sum_q (y_{hq} - \tilde{\eta}_h) = 0$$

• c. 1 point The sum of squared errors associated with this least squares estimator is the unconstrained sum of squared errors  $SSE_u$ . How would you set up a regression with dummy variables which would give you this  $SSE_u$ ?

ANSWER. The unconstrained model should be regressed in the form  $y_{mq} = \eta_m + \varepsilon_{mq}$ . String out the matrix  $\mathbf{Y}$  as a vector and for each column of  $\mathbf{Y}$  introduce a dummy variable which is 1 if the given observation was originally in this column.

• d. 2 points Next turn to the constrained model (31.2.26). If  $\mathbf{X}$  has full column rank, then it is fully identified. Writing  $\tilde{\beta}_j$  for your estimates of  $\beta_j$ , give a formula for the sum of squared errors of this estimate. By taking the first order conditions show that the estimate  $\hat{\boldsymbol{\beta}}$  is the same as the estimate in the model without replicated observations

$$(31.2.30) \quad \mathbf{z}_m = \sum_{j=1}^k x_{mj} \beta_j + \varepsilon_m,$$

where  $z_m = \bar{y}_m$  as defined above.

• e. 2 points If  $SSE_c$  is the SSE in the constrained model (31.2.26) and  $SSE_b$  the SSE in (31.2.30), show that  $SSE_c = r \cdot SSE_b + SSE_u$ .

ANSWER. For every  $m$  we have  $\sum_q (y_{mq} - \mathbf{x}_m^\top \hat{\boldsymbol{\beta}})^2 = \sum_q (y_{mq} - \bar{y}_m)^2 + r(y_{m\cdot} - \mathbf{x}_m^\top \hat{\boldsymbol{\beta}})^2$ ; therefore  $SSE_c = \sum_{m,q} (y_{mq} - \bar{y}_m)^2 + r \sum_m (y_{m\cdot} - \mathbf{x}_m^\top \hat{\boldsymbol{\beta}})^2$ ;  $\square$

• f. 3 points Write down the formula of the F-test in terms of  $SSE_u$  and  $SSE_c$  with a correct accounting of the degrees of freedom, and give this formula also in terms of  $SSE_u$  and  $SSE_b$ .

ANSWER. Unconstrained model has  $n$  parameters, and constrained model has  $k$  parameters; the number of additional “constraints” is therefore  $n - k$ . This gives the F-statistic

$$(31.2.31) \quad \frac{(SSE_c - SSE_u)/(n - k)}{SSE_u/n(r - 1)} = \frac{rSSE_b/(n - k)}{SSE_u/n(r - 1)}$$

$\square$

### 31.3. The F-Test Statistic is a Function of the Likelihood Ratio

PROBLEM 355. The critical region of the generalized likelihood ratio test can be written as

$$(31.3.1) \quad C = \{y_1, \dots, y_n : \frac{\sup_{\boldsymbol{\theta} \in \Omega} \ell(y_1, \dots, y_n; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)}{\sup_{\boldsymbol{\theta} \in \omega} \ell(y_1, \dots, y_n; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)} \geq k\},$$

where  $\omega$  refers to the null and  $\Omega$  to the alternative hypothesis (it is assumed that the hypotheses are nested, i.e.,  $\omega \subset \Omega$ ). In other words, one rejects the hypothesis if the maximal achievable likelihood level with the restriction imposed is much lower than that without the restriction. If  $\hat{\boldsymbol{\theta}}$  is the unrestricted and  $\hat{\hat{\boldsymbol{\theta}}}$  the restricted maximum likelihood estimator, then the test statistic is

$$(31.3.2) \quad LR = 2(\log \ell(\mathbf{y}, \hat{\boldsymbol{\theta}}) - \log \ell(\mathbf{y}, \hat{\hat{\boldsymbol{\theta}}})) \rightarrow \chi_i^2$$

where  $i$  is the number of restrictions. In this exercise we are proving that the F-test in the linear model is equivalent to the generalized likelihood ratio test. (You should assume here that both  $\boldsymbol{\beta}$  and  $\sigma^2$  are unknown.) All this is in [Gre97, p. 304].

• a. 1 point Since we only have constraints on  $\boldsymbol{\beta}$  and not on  $\sigma^2$ , it makes sense to first compute the concentrated likelihood function with  $\sigma^2$  concentrated out. Derive the formula for this concentrated likelihood function which is given in [Gre97, just above (6.88)].

ANSWER.

$$(31.3.3) \quad \text{Concentrated } \log \ell(\mathbf{y}; \boldsymbol{\beta}) = -\frac{n}{2} \left( 1 + \log 2\pi + \log \frac{1}{n} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right)$$

$\square$

• b. 2 points In the case of a linear restriction, show that LR is connected to the F-statistic F as follows:

$$(31.3.4) \quad LR = n \log \left( 1 + \frac{i}{n - k} F \right)$$

ANSWER.  $LR = -n \left( \log \frac{1}{n} \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} - \log \frac{1}{n} \hat{\hat{\boldsymbol{\varepsilon}}}^\top \hat{\hat{\boldsymbol{\varepsilon}}} \right) = n \log \frac{\hat{\hat{\boldsymbol{\varepsilon}}}^\top \hat{\hat{\boldsymbol{\varepsilon}}}}{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}$  In order to connect this with the statistic note that

$$(31.3.5) \quad F = \frac{n - k}{i} \left( \frac{\hat{\hat{\boldsymbol{\varepsilon}}}^\top \hat{\hat{\boldsymbol{\varepsilon}}}}{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}} - 1 \right)$$

### 31.4. Tests of Nonlinear Hypotheses

Make linear approximation, need Jacobian for this. Here is an example where nonlinear hypothesis arises naturally:

PROBLEM 356. [Gre97, Example 7.14 on p. 361]: The model

$$(31.4.1) \quad C_t = \alpha + \beta Y_t + \gamma C_{t-1} + \varepsilon_t$$

has different long run and short run propensities to consume. Give formulas for both.

ANSWER. Short-run is  $\beta$ ; to compute the long run propensity, which would prevail in stationary state when  $C_t = C_{t-1}$ , write  $C_\infty = \alpha + \beta Y_\infty + \gamma C_\infty + \varepsilon_\infty$  or  $C_\infty(1 - \gamma) = \alpha + \beta Y_\infty + \varepsilon_\infty$  or  $C_\infty = \alpha/(1 - \gamma) + \beta/(1 - \gamma) Y_\infty + \varepsilon_t/(1 - \gamma)$ . Therefore long run propensity is  $\delta = \beta/(1 - \gamma)$ .

### 31.5. Choosing Between Nonnested Models

Throwing all regressors into the same regression is a straightforward way out but not very good. J-test (the J comes from “joint”) is better: throw the predicted value of one of the two models as a regressor into the other model and test whether the predicted value has a nonzero coefficient. Here is more detail: if the null hypothesis is that model 1 is right, then throw the predicted value of model 2 into model 1 and test the null hypothesis that the coefficient of this predicted value is zero. If Model 1 is right, then this additional regressor leaves all other estimators unbiased and the true coefficient of the additional regressor is 0. If Model 2 is right, then asymptotically, this additional regressor should be the only regressor in the combined model with a nonzero coefficient (its coefficient is = 1 asymptotically, and all the other regressors should have coefficient zero.) Whenever nonnested hypotheses are tested, it is possible that both hypotheses are rejected, or that neither hypothesis is rejected by this criterion.



## CHAPTER 32

## Instrumental Variables

Compare here [DM93, chapter 7] and [Gre97, Section 6.7.8]. Greene first introduces the simple instrumental variables estimator and then shows that the generalized one picks out the best linear combinations for forming simple instruments. I will follow [DM93] and first introduce the generalized instrumental variables estimator, and then go down to the simple one.

In this chapter, we will discuss a sequence of models  $\mathbf{y}_n = \mathbf{X}_n\boldsymbol{\beta} + \boldsymbol{\varepsilon}_n$ , where  $\boldsymbol{\varepsilon}_n \sim (\mathbf{o}_n, \sigma^2 \mathbf{I}_n)$ , and  $\mathbf{X}_n$  are  $n \times k$ -matrices of random regressors, and the number of observations  $n \rightarrow \infty$ . We do not make the assumption  $\text{plim } \frac{1}{n} \mathbf{X}_n^\top \boldsymbol{\varepsilon}_n = \mathbf{o}$  which would ensure consistency of the OLS estimator (compare Problem 328). Instead, a sequence of  $n \times m$  matrices of (random or nonrandom) “instrumental variables”  $\mathbf{W}_n$  is available which satisfies the following three conditions:

$$(32.0.1) \quad \text{plim } \frac{1}{n} \mathbf{W}_n^\top \boldsymbol{\varepsilon}_n = \mathbf{o}$$

$$(32.0.2) \quad \text{plim } \frac{1}{n} \mathbf{W}_n^\top \mathbf{W}_n = \mathbf{Q} \quad \text{exists, is nonrandom and nonsingular}$$

$$(32.0.3) \quad \text{plim } \frac{1}{n} \mathbf{W}_n^\top \mathbf{X}_n = \mathbf{D} \quad \text{exists, is nonrandom and has full column rank}$$

Full column rank in (32.0.3) is only possible if  $m \geq k$ .

In this situation, regression of  $\mathbf{y}$  on  $\mathbf{X}$  is inconsistent. But if one regresses  $\mathbf{y}$  on the projection of  $\mathbf{X}$  on  $\text{R}[\mathbf{W}]$ , the column space of  $\mathbf{W}$ , one obtains a consistent estimator. This is called the instrumental variables estimator.

If  $\mathbf{x}_i$  is the  $i$ th column vector of  $\mathbf{X}$ , then  $\mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{x}_i$  is the projection of  $\mathbf{x}_i$  on the space spanned by the columns of  $\mathbf{W}$ . Therefore the matrix  $\mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X}$  consists of the columns of  $\mathbf{X}$  projected on  $\text{R}[\mathbf{W}]$ . This is what we meant by the projection of  $\mathbf{X}$  on  $\text{R}[\mathbf{W}]$ . With these projections as regressors, the vector of regression coefficients becomes the “generalized instrumental variables estimator”

$$(32.0.4) \quad \tilde{\boldsymbol{\beta}} = \left( \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}$$

PROBLEM 357. 3 points We are in the model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  and we have a matrix  $\mathbf{W}$  of “instrumental variables” which satisfies the following three conditions:

$\text{plim } \frac{1}{n} \mathbf{W}^\top \boldsymbol{\varepsilon} = \mathbf{o}$ ,  $\text{plim } \frac{1}{n} \mathbf{W}^\top \mathbf{W} = \mathbf{Q}$  exists, is nonrandom and positive definite, and  $\text{plim } \frac{1}{n} \mathbf{W}^\top \mathbf{X} = \mathbf{D}$  exists, is nonrandom and has full column rank. Show that instrumental variables estimator

$$(32.0.5) \quad \tilde{\boldsymbol{\beta}} = \left( \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}$$

is consistent. Hint: Write  $\tilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta} = \mathbf{B}_n \cdot \frac{1}{n} \mathbf{W}^\top \boldsymbol{\varepsilon}$  and show that the sequence matrices  $\mathbf{B}_n$  has a plim.

ANSWER. Write it as

$$\begin{aligned} \tilde{\boldsymbol{\beta}}_n &= \left( \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) \\ &= \boldsymbol{\beta} + \left( \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \boldsymbol{\varepsilon} \\ &= \boldsymbol{\beta} + \left( \left( \frac{1}{n} \mathbf{X}^\top \mathbf{W} \right) \left( \frac{1}{n} \mathbf{W}^\top \mathbf{W} \right)^{-1} \left( \frac{1}{n} \mathbf{W}^\top \mathbf{X} \right) \right)^{-1} \left( \frac{1}{n} \mathbf{X}^\top \mathbf{W} \right) \left( \frac{1}{n} \mathbf{W}^\top \mathbf{W} \right)^{-1} \frac{1}{n} \mathbf{W}^\top \boldsymbol{\varepsilon}, \end{aligned}$$

i.e., the  $\mathbf{B}_n$  and  $\mathbf{B}$  of the hint are as follows:

$$\begin{aligned} \mathbf{B}_n &= \left( \left( \frac{1}{n} \mathbf{X}^\top \mathbf{W} \right) \left( \frac{1}{n} \mathbf{W}^\top \mathbf{W} \right)^{-1} \left( \frac{1}{n} \mathbf{W}^\top \mathbf{X} \right) \right)^{-1} \left( \frac{1}{n} \mathbf{X}^\top \mathbf{W} \right) \left( \frac{1}{n} \mathbf{W}^\top \mathbf{W} \right)^{-1} \\ \mathbf{B} &= \text{plim } \mathbf{B}_n = (\mathbf{D}^\top \mathbf{Q}^{-1} \mathbf{D})^{-1} \mathbf{D}^\top \mathbf{Q}^{-1} \end{aligned}$$

PROBLEM 358. Assume  $\text{plim } \frac{1}{n} \mathbf{X}^\top \mathbf{X}$  exists, and  $\text{plim } \frac{1}{n} \mathbf{X}^\top \boldsymbol{\varepsilon}$  exists. (We do not need the existence, not that the first is nonsingular and the second zero). Show that  $\sigma^2$  can be estimated consistently by  $s^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})$ .

ANSWER.  $\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} - \mathbf{X}\tilde{\boldsymbol{\beta}} = \boldsymbol{\varepsilon} - \mathbf{X}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})$ . Therefore

$$\frac{1}{n} (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) = \frac{1}{n} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon} - \frac{2}{n} \boldsymbol{\varepsilon}^\top \mathbf{X}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}) + (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \left( \frac{1}{n} \mathbf{X}^\top \mathbf{X} \right) (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}).$$

All summands have plims, the plim of the first is  $\sigma^2$  and those of the other two are zero.

PROBLEM 359. In the situation of Problem 357, add the stronger assumption  $\frac{1}{n} \mathbf{W}^\top \boldsymbol{\varepsilon} \rightarrow \text{N}(\mathbf{o}, \sigma^2 \mathbf{Q})$ , and show that  $\sqrt{n}(\tilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \rightarrow \text{N}(\mathbf{o}, \sigma^2 (\mathbf{D}^\top \mathbf{Q}^{-1} \mathbf{D})^{-1})$

ANSWER.  $\tilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta} = \mathbf{B}_n \frac{1}{n} \mathbf{W}_n^\top \boldsymbol{\varepsilon}_n$ , therefore  $\sqrt{n}(\tilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) = \mathbf{B}_n n^{-1/2} \mathbf{W}_n^\top \boldsymbol{\varepsilon}_n \rightarrow \text{BN}(\mathbf{o}, \sigma^2 \mathbf{Q} \mathbf{B} \mathbf{B}^\top)$ . Since  $\mathbf{B} = (\mathbf{D}^\top \mathbf{Q}^{-1} \mathbf{D})^{-1} \mathbf{D}^\top \mathbf{Q}^{-1}$ , the result follows.

From Problem 359 follows that for finite samples approximately  $\tilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta} \sim \text{N}(\mathbf{o}, \frac{\sigma^2}{n} (\mathbf{D}^\top \mathbf{Q}^{-1} \mathbf{D})^{-1})$ . Since  $\frac{1}{n} (\mathbf{D}^\top \mathbf{Q}^{-1} \mathbf{D})^{-1} = (n \mathbf{D}^\top (n \mathbf{Q})^{-1} n \mathbf{D})^{-1}$ ,  $\text{MSE}[\tilde{\boldsymbol{\beta}}_n]$  can be estimated by  $s^2 \left( \mathbf{X}^\top \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X} \right)^{-1}$

The estimator (32.0.4) is sometimes called the two stages least squares estimator because the projection of  $\mathbf{X}$  on the column space of  $\mathbf{W}$  can be considered the predicted

values if one regresses every column of  $\mathbf{X}$  on  $\mathbf{W}$ . I.e., instead of regressing  $\mathbf{y}$  on  $\mathbf{X}$  one regresses  $\mathbf{y}$  on those linear combinations of the columns of  $\mathbf{W}$  which best approximate the columns of  $\mathbf{X}$ . Here is more detail: the matrix of estimated coefficients in the first regression is  $\hat{\mathbf{\Pi}} = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X}$ , and the predicted values in this regression are  $\hat{\mathbf{X}} = \mathbf{W} \hat{\mathbf{\Pi}} = \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X}$ . The second regression, which regresses  $\mathbf{y}$  on  $\hat{\mathbf{X}}$ , gives the coefficient vector

$$(32.0.6) \quad \tilde{\boldsymbol{\beta}} = (\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top \mathbf{y}.$$

If you plug this in you see this is exactly (32.0.4) again.

Now let's look at the geometry of instrumental variable regression of one variable  $\mathbf{y}$  on one other variable  $\mathbf{x}$  with  $\mathbf{w}$  as an instrument. The specification is  $\mathbf{y} = \mathbf{x}\beta + \boldsymbol{\varepsilon}$ . On p. 347 we visualized the asymptotic results if  $\boldsymbol{\varepsilon}$  is asymptotically orthogonal to  $\mathbf{x}$ . Now let us assume  $\boldsymbol{\varepsilon}$  is asymptotically not orthogonal to  $\mathbf{x}$ . One can visualize this as three vectors, again normalized by dividing by  $\sqrt{n}$ , but now even in the asymptotic case the  $\boldsymbol{\varepsilon}$ -vector is not orthogonal to  $\mathbf{x}$ . (Draw  $\boldsymbol{\varepsilon}$  vertically, and make  $\mathbf{x}$  long enough that  $\beta < 1$ .) We assume  $n$  is large enough so that the asymptotic results hold for the sample already (or, perhaps better, that the difference between the sample and its plim is only infinitesimal). Therefore the OLS regression, with estimates  $\beta$  by  $\mathbf{x}^\top \mathbf{y} / \mathbf{x}^\top \mathbf{x}$ , is inconsistent. Let  $O$  be the origin,  $A$  the point on the  $\mathbf{x}$ -vector where  $\boldsymbol{\varepsilon}$  branches off (i.e., the end of  $\mathbf{x}\beta$ ), furthermore let  $B$  be the point on the  $\mathbf{x}$ -vector where the orthogonal projection of  $\mathbf{y}$  comes down, and  $C$  the end of the  $\mathbf{x}$ -vector. Then  $\mathbf{x}^\top \mathbf{y} = \overline{OC} \overline{OB}$  and  $\mathbf{x}^\top \mathbf{x} = \overline{OC}^2$ , therefore  $\mathbf{x}^\top \mathbf{y} / \mathbf{x}^\top \mathbf{x} = \overline{OB} / \overline{OC}$ , which would be the  $\beta$  if the errors were orthogonal. Now introduce a new variable  $\mathbf{w}$  which is orthogonal to the errors. (Since  $\boldsymbol{\varepsilon}$  is vertical,  $\mathbf{w}$  is on the horizontal axis.) Call  $D$  the projection of  $\mathbf{y}$  on  $\mathbf{w}$ , which is the prolongation of the vector  $\boldsymbol{\varepsilon}$ , and call  $E$  the end of the  $\mathbf{w}$ -vector, and call  $F$  the projection of  $\mathbf{x}$  on  $\mathbf{w}$ . Then  $\mathbf{w}^\top \mathbf{y} = \overline{OE} \overline{OD}$ , and  $\mathbf{w}^\top \mathbf{x} = \overline{OE} \overline{OF}$ . Therefore  $\mathbf{w}^\top \mathbf{y} / \mathbf{w}^\top \mathbf{x} = (\overline{OE} \overline{OD}) / (\overline{OE} \overline{OF}) = \overline{OD} / \overline{OF} = \overline{OA} / \overline{OC} = \beta$ . Or geometrically it is obvious that the regression of  $\mathbf{y}$  on the projection of  $\mathbf{x}$  on  $\mathbf{w}$  will give the right  $\hat{\beta}$ . One also sees here why the  $s^2$  based on this second regression is inconsistent.

If I allow two instruments, the two instruments must be in the horizontal plane perpendicular to the vector  $\boldsymbol{\varepsilon}$  which is assumed still vertical. Here we project  $\mathbf{x}$  on this horizontal plane and then regress the  $\mathbf{y}$ , which stays where it is, on this  $\mathbf{x}$ . In this way the residuals have the right direction!

What if there is one instrument, but it does not lie in the same plane as  $\mathbf{x}$  and  $\mathbf{y}$ ? This is the most general case as long as there is only one regressor and one instrument. This instrument  $\mathbf{w}$  must lie somewhere in the horizontal plane. We have to project  $\mathbf{x}$  on it, and then regress  $\mathbf{y}$  on this projection. Look at it this way: take the plane orthogonal to  $\mathbf{w}$  which goes through point  $C$ . The projection of  $\mathbf{x}$  on  $\mathbf{w}$  is the intersection of the ray generated by  $\mathbf{w}$  with this plane. Now move this

plane parallel until it intersects point  $A$ . Then the intersection with the  $\mathbf{w}$ -ray is the projection of  $\mathbf{y}$  on  $\mathbf{w}$ . But this latter plane contains  $\boldsymbol{\varepsilon}$ , since  $\boldsymbol{\varepsilon}$  is orthogonal to  $\mathbf{w}$ . This makes sure that the regression gives the right results.

PROBLEM 360. 4 points The asymptotic MSE matrix of the instrumental variables estimator with  $\mathbf{W}$  as matrix of instruments is  $\sigma^2 \text{plim}(\mathbf{X}^\top \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X})^{-1}$ . Show that if one adds more instruments, then this asymptotic MSE-matrix can only decrease. It is sufficient to show that the inequality holds before going over to plim, i.e., if  $\mathbf{W} = [\mathbf{U} \quad \mathbf{V}]$ , then

$$(32.0.7) \quad (\mathbf{X}^\top \mathbf{U} (\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{X})^{-1} - (\mathbf{X}^\top \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X})^{-1}$$

is nonnegative definite. Hints: (1) Use theorem A.5.5 in the Appendix (proof not required). (2) Note that  $\mathbf{U} = \mathbf{W}\mathbf{G}$  for some  $\mathbf{G}$ . Can you write this  $\mathbf{G}$  in partitioned matrix form? (3) Show that, whatever  $\mathbf{W}$  and  $\mathbf{G}$ ,  $\mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{W} \mathbf{G} (\mathbf{G}^\top \mathbf{W}^\top \mathbf{W} \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{W}^\top$  is idempotent.

ANSWER.

$$(32.0.8) \quad \mathbf{u} = [\mathbf{u} \quad \mathbf{v}] \begin{bmatrix} \mathbf{I} \\ \mathbf{O} \end{bmatrix} = \mathbf{W}\mathbf{G} \quad \text{where} \quad \mathbf{G} = \begin{bmatrix} \mathbf{I} \\ \mathbf{O} \end{bmatrix}.$$

PROBLEM 361. 2 points Show: if a matrix  $\mathbf{D}$  has full column rank and is square then it has an inverse.

ANSWER. Here you need that column rank is row rank: if  $\mathbf{D}$  has full column rank it also has full row rank. And to make the proof complete you need: if  $\mathbf{A}$  has a left inverse  $\mathbf{L}$  and a right inverse  $\mathbf{R}$ , then  $\mathbf{L}$  is the only left inverse and  $\mathbf{R}$  the only right inverse and  $\mathbf{L} = \mathbf{R}$ . Proof:  $\mathbf{L} = \mathbf{L}(\mathbf{A}\mathbf{R}) = (\mathbf{L}\mathbf{A})\mathbf{R} = \mathbf{R}$ .

PROBLEM 362. 2 points If  $\mathbf{W}^\top \mathbf{X}$  is square and has full column rank, then it is nonsingular. Show that in this case (32.0.4) simplifies to the "simple" instrumental variables estimator:

$$(32.0.9) \quad \tilde{\boldsymbol{\beta}} = (\mathbf{W}^\top \mathbf{X})^{-1} \mathbf{W}^\top \mathbf{y}$$

ANSWER. In this case the big inverse can be split into three:

$$(32.0.10) \quad \tilde{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y} =$$

$$(32.0.11) \quad = (\mathbf{W}^\top \mathbf{X})^{-1} \mathbf{W}^\top \mathbf{W} (\mathbf{X}^\top \mathbf{W})^{-1} \mathbf{X}^\top \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}$$

PROBLEM 363. We only have one regressor with intercept, i.e.,  $\mathbf{X} = [\mathbf{1} \ \mathbf{x}]$ , and we have one instrument  $\mathbf{w}$  for  $\mathbf{x}$  (while the constant term is its own instrument), i.e.,  $\mathbf{W} = [\mathbf{1} \ \mathbf{w}]$ . Show that the instrumental variables estimators for slope and intercept are

$$(32.0.12) \quad \tilde{\beta} = \frac{\sum (w_t - \bar{w})(y_t - \bar{y})}{\sum (w_t - \bar{w})(x_t - \bar{x})}$$

$$(32.0.13) \quad \tilde{\alpha} = \bar{y} - \tilde{\beta}\bar{x}$$

*Hint: the math is identical to that in question 200.*

PROBLEM 364. 2 points Show that, if there are as many instruments as there are observations, then the instrumental variables estimator (32.0.4) becomes identical to OLS.

ANSWER. In this case  $\mathbf{W}$  has an inverse, therefore the projection on  $\mathbf{R}[\mathbf{W}]$  is the identity. Staying in the algebraic paradigm,  $(\mathbf{W}^\top \mathbf{W})^{-1} = \mathbf{W}^{-1}(\mathbf{W}^\top)^{-1}$ .  $\square$

An implication of Problem 364 is that one must be careful not to include too many instruments if one has a small sample. Asymptotically it is better to have more instruments, but for  $n = m$ , the instrumental variables estimator is equal to OLS, i.e., the sequence of instrumental variables estimators starts at the (inconsistent) OLS. If one uses fewer instruments, then the asymptotic  $\mathcal{MSE}$  matrix is not so good, but one may get a sequence of estimators which moves away from the inconsistent OLS more quickly.

## APPENDIX A

## Matrix Formulas

In this Appendix, efforts are made to give some of the familiar matrix lemmas in their most general form. The reader should be warned: the concept of a deficiency matrix and the notation which uses a thick fraction line multiplication with a scalar g-inverse are my own.

## A.1. A Fundamental Matrix Decomposition

**THEOREM A.1.1.** *Every matrix  $B$  which is not the null matrix can be written as a product of two matrices  $B = CD$ , where  $C$  has a left inverse  $L$  and  $D$  a right inverse  $R$ , i.e.,  $LC = DR = I$ . This identity matrix is  $r \times r$ , where  $r$  is the rank of  $B$ .*

A proof is in [Rao73, p. 19]. This is the fundamental theorem of algebra, that every homomorphism can be written as a product of epimorphism and monomorphism, together with the fact that all epimorphisms and monomorphisms split, i.e., have one-sided inverses.

One such factorization is given by the singular value theorem: If  $B = P^T \Lambda Q$  is the svd as in Theorem A.9.2, then one might set e.g.  $C = P^T \Lambda$  and  $D = Q$ , consequently  $L = \Lambda^{-1} P$  and  $R = Q^T$ . In this decomposition, the first row/column carries the largest weight and gives the best approximation in a least squares sense, etc.

The trace of a square matrix is defined as the sum of its diagonal elements. The rank of a matrix is defined as the number of its linearly independent rows, which is equal to the number of its linearly independent columns (row rank = column rank).

**THEOREM A.1.2.**  $\text{tr } BC = \text{tr } CB$ .

**PROBLEM 365.** *Prove theorem A.1.2.*

**PROBLEM 366.** *Use theorem A.1.1 to prove that if  $BB = B$ , then  $\text{rank } B = \text{tr } B$ .*

**ANSWER.** Premultiply the equation  $CD = CDCD$  by  $L$  and postmultiply it by  $R$  to get  $DC = I_r$ . This is useful for the trace:  $\text{tr } B = \text{tr } CD = \text{tr } DC = \text{tr } I_r = r$ . I have this proof from [Rao73, p. 28].  $\square$

**THEOREM A.1.3.**  $B = O$  if and only if  $B^T B = O$ .

## A.2. The Spectral Norm of a Matrix

The spectral norm of a matrix extends the Euclidean norm  $\|z\|$  from vectors to matrices. Its definition is  $\|A\| = \max_{\|z\|=1} \|Az\|$ . This spectral norm is the maximum singular value  $\mu_{max}$ , and if  $A$  is square, then  $\|A^{-1}\| = 1/\mu_{min}$ . It is a true norm, i.e.,  $\|A\| = 0$  if and only if  $A = O$ , furthermore  $\|\lambda A\| = |\lambda| \cdot \|A\|$ , and the triangle inequality  $\|A + B\| \leq \|A\| + \|B\|$ . In addition, it obeys  $\|AB\| \leq \|A\| \cdot \|B\|$ .

**PROBLEM 367.** *Show that the spectral norm is the maximum singular value.*

**ANSWER.** Use the definition

$$(A.2.1) \quad \|A\|^2 = \max \frac{z^T A^T A z}{z^T z}.$$

Write  $A = P^T \Lambda Q$  as in (A.9.1). Then  $z^T A^T A z = z^T Q^T \Lambda^2 Q z$ . Therefore we can first search for a  $z$  in the form  $z = Q^T x$  which attains this maximum. Proof: for every  $z$  which has a nonzero value in the numerator of (A.2.1), set  $x = Qz$ . Then  $x \neq o$ , and  $Q^T x$  attains the same value as  $z$  in the numerator of (A.2.1), and a smaller or equal value in the denominator. Therefore one can restrict the search for the maximum argument to vectors of the form  $Q^T x$ . But for this the objective function becomes  $\frac{x^T \Lambda^2 x}{x^T x}$ , which is maximized by  $x = i_1$ , the first unit vector (the  $i$ th column vector of the unit matrix). Therefore the squared spectral norm is  $\lambda_{ii}^2$ , and therefore the spectral norm itself is  $\lambda_{ii}$ .

## A.3. Inverses and g-Inverses of Matrices

A g-inverse of a matrix  $A$  is any matrix  $A^-$  satisfying

$$(A.3.1) \quad A = AA^-A.$$

It always exists but is not always unique. If  $A$  is square and nonsingular, then  $A^{-1}$  is its only g-inverse.

**PROBLEM 368.** *Show that a symmetric matrix  $\Omega$  has a g-inverse which is also symmetric.*

**ANSWER.** Use  $\Omega^- \Omega \Omega^{-T}$ .

The definition of a g-inverse is apparently due to [Rao62]. It is sometimes called the “conditional inverse” [Gra83, p. 129]. This g-inverse, and not the Moore-Penrose generalized inverse or pseudoinverse  $A^+$ , is needed for the linear model. The Moore-Penrose generalized inverse is a g-inverse that in addition satisfies  $A^+ A A^+ = A^+$  and  $A A^+$  as well as  $A^+ A$  symmetric. It always exists and is also unique, but the additional requirements are burdensome ballast. [Gre97, pp. 44-5] also advocates the Moore-Penrose inverse, but he does not really use it. If he were to try to use it, he would probably soon discover that it is not appropriate. The book [Alb72] de-

the linear model with the Moore-Penrose inverse. It is a good demonstration of how complicated everything gets if one uses an inappropriate mathematical tool.

PROBLEM 369. Use theorem A.1.1 to prove that every matrix has a g-inverse.

ANSWER. Simple: a null matrix has its transpose as g-inverse, and if  $A \neq O$  then  $RL$  is such a g-inverse.  $\square$

The g-inverse of a number is its inverse if the number is nonzero, and is arbitrary otherwise. Scalar expressions written as fractions are in many cases the multiplication by a g-inverse. We will use a fraction with a thick horizontal rule to indicate where this is the case. In other words, by definition,

$$(A.3.2) \quad \frac{a}{b} = b^{-}a. \quad \text{Compare that with the ordinary fraction } \frac{a}{b}.$$

This idiosyncratic notation allows to write certain theorems in a more concise form, but it requires more work in the proofs, because one has to consider the additional case that the denominator is zero. Theorems A.5.8 and A.8.2 are examples.

THEOREM A.3.1. If  $B = AA^{-}B$  holds for one g-inverse  $A^{-}$  of  $A$ , then it holds for all g-inverses. If  $A$  is symmetric and  $B = AA^{-}B$ , then also  $B^{\top} = B^{\top}A^{-}A$ . If  $B = BA^{-}A$  and  $C = AA^{-}C$  then  $BA^{-}C$  is independent of the choice of g-inverses.

PROOF. Assume the identity  $B = AA^{+}B$  holds for some fixed g-inverse  $A^{+}$  (which may be, as the notation suggests, the Moore Penrose g-inverse, but this is not necessary), and let  $A^{-}$  be an different g-inverse. Then  $AA^{-}B = AA^{-}AA^{+}B = AA^{+}B = B$ . For the second statement one merely has to take transposes and note that a matrix is a g-inverse of a symmetric  $A$  if and only if its transpose is. For the third statement:  $BA^{+}C = BA^{-}AA^{+}AA^{-}C = BA^{-}AA^{-}C = BA^{-}C$ . Here  $+$  signifies a different g-inverse; again, it is not necessarily the Moore-Penrose one.  $\square$

PROBLEM 370. Show that  $x$  satisfies  $x = Ba$  for some  $a$  if and only if  $x = BB^{-}x$ .

THEOREM A.3.2. Both  $A^{\top}(AA^{\top})^{-}$  and  $(A^{\top}A)^{-}A$  are g-inverses of  $A$ .

PROOF. We have to show

$$(A.3.3) \quad A = AA^{\top}(AA^{\top})^{-}A$$

which is [Rao73, (1b.5.5) on p. 26]. Define  $D = A - AA^{\top}(AA^{\top})^{-}A$  and show, by multiplying out, that  $DD^{\top} = O$ .  $\square$

### A.4. Deficiency Matrices

Here is again some idiosyncratic terminology and notation. It gives an explicit algebraic formulation for something that is often done implicitly or in a geometric paradigm. A matrix  $G$  will be called a “left deficiency matrix” of  $S$ , in symbols  $G \perp S$ , if  $GS = O$ , and for all  $Q$  with  $QS = O$  there is an  $X$  with  $Q = XG$ . This factorization property is an algebraic formulation of the geometric concept of a null space. It is symmetric in the sense that  $G \perp S$  is also equivalent with:  $GS = O$  and for all  $R$  with  $GR = O$  there is a  $Y$  with  $R = SY$ . In other words,  $G \perp S$  and  $S^{\top} \perp G^{\top}$  are equivalent.

This symmetry follows from the following characterization of a deficiency matrix which is symmetric:

THEOREM A.4.1.  $T \perp U$  iff  $TU = O$  and  $T^{\top}T + UU^{\top}$  nonsingular.

PROOF. This proof here seems terribly complicated. There must be a simpler way. Proof of “ $\Rightarrow$ ”: Assume  $T \perp U$ . Take any  $\gamma$  with  $\gamma^{\top}T^{\top}T\gamma + \gamma^{\top}UU^{\top}\gamma = 0$ , i.e.,  $T\gamma = o$  and  $\gamma^{\top}U = o^{\top}$ . From this one can show that  $\gamma = o$ : since  $T\gamma = o$ , there is a  $\xi$  with  $\gamma = U\xi$ , therefore  $\gamma^{\top}\gamma = \gamma^{\top}U\xi = 0$ . To prove “ $\Leftarrow$ ” assume  $TU = O$  and  $T^{\top}T + UU^{\top}$  is nonsingular. To show that  $T \perp U$  take any  $B$  with  $BU = O$ . Then  $B = B(T^{\top}T + UU^{\top})(T^{\top}T + UU^{\top})^{-1} = BT^{\top}T(T^{\top}T + UU^{\top})^{-1}$ . In the same way one gets  $T = TT^{\top}T(T^{\top}T + UU^{\top})^{-1}$ . Premultiply this last equation by  $T^{\top}T(T^{\top}TT^{\top}T) - T^{\top}$  and use theorem A.3.2 to get  $T^{\top}T(T^{\top}TT^{\top}T) - T^{\top}T = T^{\top}T(T^{\top}T + UU^{\top})^{-1}$ . Inserting this into the equation for  $B$  gives  $B = BT^{\top}T(T^{\top}TT^{\top}T) - T^{\top}T$ , i.e.,  $B$  factors over  $T$ .

The R/Spplus-function Null gives the transpose of a deficiency matrix.

THEOREM A.4.2. If for all  $Y$ ,  $BY = O$  implies  $AY = O$ , then a  $X$  exists with  $A = XB$ .

PROBLEM 371. Prove theorem A.4.2.

ANSWER. Let  $B \perp C$ . Choosing  $Y = B$  follows  $AB = O$ , hence  $X$  exists.

PROBLEM 372. Show that  $I - SS^{-} \perp S$ .

ANSWER. Clearly,  $(I - SS^{-})S = O$ . Now if  $QS = O$ , then  $Q = Q(I - SS^{-})$ , i.e., the whose existence is postulated in the definition of a deficiency matrix is  $Q$  itself.

PROBLEM 373. Show that  $S \perp U$  if and only if  $S$  is a matrix with maximal rank which satisfies  $SU = O$ . In other words, one cannot add linearly independent rows to  $S$  in such a way that the new matrix still satisfies  $TU = O$ .

ANSWER. First assume  $S \perp U$  and take any additional row  $t^\top$  so that  $\begin{bmatrix} S \\ t^\top \end{bmatrix} U = \begin{bmatrix} O \\ o^\top \end{bmatrix}$ . Then exists a  $\begin{bmatrix} Q \\ r \end{bmatrix}$  such that  $\begin{bmatrix} S \\ t^\top \end{bmatrix} = \begin{bmatrix} Q \\ r \end{bmatrix} S$ , i.e.,  $SQ = S$ , and  $t^\top = r^\top S$ . But this last equation means that  $t^\top$  is a linear combination of the rows of  $S$  with the  $r_i$  as coefficients. Now conversely, assume  $S$  is such that one cannot add a linearly independent row  $t^\top$  such that  $\begin{bmatrix} S \\ t^\top \end{bmatrix} U = \begin{bmatrix} O \\ o^\top \end{bmatrix}$ , and let  $PU = O$ . Then all rows of  $P$  must be linear combinations of rows of  $S$  (otherwise one could add such a row to  $S$  and get the result which was just ruled out), therefore  $P = SS$  where  $A$  is the matrix of coefficients of these linear combinations.  $\square$

The deficiency matrix is not unique, but we will use the concept of a deficiency matrix in a formula only then when this formula remains correct for every deficiency matrix. One can make deficiency matrices unique if one requires them to be projection matrices.

PROBLEM 374. Given  $X$  and a symmetric nonnegative definite  $\Omega$  such that  $X = \Omega W$  for some  $W$ . Show that  $X \perp U$  if and only if  $X^\top \Omega^- X \perp U$ .

ANSWER. One has to show that  $XY = O$  is equivalent to  $X^\top \Omega^- XY = O$ .  $\Rightarrow$  clear; for  $\Leftarrow$  note that  $X^\top \Omega^- X = W^\top \Omega W$ , therefore  $XY = \Omega WY = \Omega W(W^\top \Omega W)^- W^\top \Omega WY = \Omega W(W^\top \Omega W)^- X^\top \Omega^- XY = O$ .  $\square$

A matrix is said to have full column rank if all its columns are linearly independent, and full row rank if its rows are linearly independent. The deficiency matrix provides a “holistic” definition for which it is not necessary to look at single rows and columns.  $X$  has full column rank if and only if  $X \perp O$ , and full row rank if and only if  $O \perp X$ .

PROBLEM 375. Show that the following three statements are equivalent: (1)  $X$  has full column rank, (2)  $X^\top X$  is nonsingular, and (3)  $X$  has a left inverse.

ANSWER. Here use  $X \perp O$  as the definition of “full column rank.” Then (1)  $\Leftrightarrow$  (2) is theorem A.4.1. Now (1)  $\Rightarrow$  (3): Since  $IO = O$ , a  $P$  exists with  $I = PX$ . And (3)  $\Rightarrow$  (1): if a  $P$  exists with  $I = PX$ , then any  $Q$  with  $QO = O$  can be factored over  $X$ , simply say  $Q = QPX$ .  $\square$

Note that the usual solution of linear matrix equations with g-inverses involves a deficiency matrix:

THEOREM A.4.3. The solution of the consistent matrix equation  $TX = A$  is

$$(A.4.1) \quad X = T^- A + UW$$

where  $T \perp U$  and  $W$  is arbitrary.

PROOF. Given consistency, i.e., the existence of at least one  $Z$  with  $TZ = A$ , (A.4.1) defines indeed a solution, since  $TX = TT^-TZ$ . Conversely, if  $Y$  satisfies  $TY = A$ , then  $T(Y - T^-A) = O$ , therefore  $Y - T^-A = UW$  for some  $W$ .

THEOREM A.4.4. Let  $L \perp T \perp U$  and  $J \perp HU \perp R$ ; then

$$\begin{bmatrix} L & O \\ -JHT^- & J \end{bmatrix} \perp \begin{bmatrix} T \\ H \end{bmatrix} \perp UR.$$

PROOF. First deficiency relation: Since  $I - TT^- = UW$  for some  $W$ ,  $-JHT^- JH = O$ , therefore the matrix product is zero. Now assume  $\begin{bmatrix} A & B \\ T & H \end{bmatrix} \begin{bmatrix} T \\ H \end{bmatrix} = \begin{bmatrix} A \\ B \end{bmatrix}$ . Then  $BHU = O$ , i.e.,  $B = DJ$  for some  $D$ . Then  $AT = -DJH$ , which has as general solution  $A = -DJHT^- + CL$  for some  $C$ . This together gives  $\begin{bmatrix} A & B \\ T & H \end{bmatrix} = \begin{bmatrix} C & D \\ -JHT^- & J \end{bmatrix} \begin{bmatrix} L & O \\ -JHT^- & J \end{bmatrix}$ . Now the second deficiency relation: clearly the product of the matrices is zero. If  $M$  satisfies  $TM = O$ , then  $M = UN$  for some  $N$ . If  $M$  furthermore satisfies  $HM = O$ , then  $HUN = O$ , therefore  $N = RP$  for some  $P$ , therefore  $M = URP$ .

THEOREM A.4.5. Assume  $\Omega$  is nonnegative definite symmetric and  $K$  is such that  $K\Omega$  is defined. Then the matrix

$$(A.4.2) \quad \Xi = \Omega - \Omega K^\top (K\Omega K^\top)^- K\Omega$$

has the following properties:

- (1)  $\Xi$  does not depend on the choice of g-inverse of  $K\Omega K^\top$  used in (A.4.2).
- (2) Any g-inverse of  $\Omega$  is also a g-inverse of  $\Xi$ , i.e.  $\Xi\Omega^- \Xi = \Xi$ .
- (3)  $\Xi$  is nonnegative definite and symmetric.
- (4) For every  $P \perp \Omega$  follows  $\begin{bmatrix} K \\ P \end{bmatrix} \perp \Xi$
- (5) If  $T$  is any other right deficiency matrix of  $\begin{bmatrix} K \\ P \end{bmatrix}$ , i.e., if  $\begin{bmatrix} K \\ P \end{bmatrix} \perp T$ , then

$$(A.4.3) \quad \Xi = T(T^\top \Omega^- T)^- T^\top.$$

Hint: show that any  $D$  satisfying  $\Xi = TDT^\top$  is a g-inverse of  $T^\top \Omega^- T$ .

In order to apply (A.4.3) show that the matrix  $T = SK$  where  $K \perp S$  and  $PS \perp K$  is a right deficiency matrix of  $\begin{bmatrix} K \\ P \end{bmatrix}$ .

Proof of theorem A.4.5: Independence of choice of g-inverse follows from theorem A.5.10. That  $\Omega^-$  is a g-inverse is also an immediate consequence of theorem A.5. From the factorization  $\Xi = \Xi\Omega^- \Xi$  follows also that  $\Xi$  is nnd symmetric (since every nnd symmetric  $\Omega$  also has a symmetric nnd g-inverse). (4) Deficiency proper

From  $\begin{bmatrix} K \\ P \end{bmatrix} Q = O$  follows  $KQ = O$  and  $PQ = O$ . From this second equation and  $P \perp \Omega$  follows  $Q = \Omega R$  for some  $R$ . Since  $K\Omega R = KQ = O$ , it follows  $Q = \Omega R = (\Omega - \Omega K^\top (K\Omega K^\top)^{-1} K\Omega)R$ .

Proof of (5): Since  $\begin{bmatrix} K \\ P \end{bmatrix} \Xi = O$  it follows  $\Xi = TA$  for some  $A$ , and therefore  $\Xi = \Xi\Omega^{-1}\Xi = T\Omega^{-1}A^\top T^\top = TDT^\top$  where  $D = \Omega^{-1}A^\top$ .

Before going on we need a lemma. Since  $(I - \Omega\Omega^{-1})\Omega = O$ , there exists a  $N$  with  $I - \Omega\Omega^{-1} = NP$ , therefore  $T - \Omega\Omega^{-1}T = NP^\top T = O$  or

$$(A.4.4) \quad T = \Omega\Omega^{-1}T$$

Using (A.4.4) one can show the hint: that any  $D$  satisfying  $\Xi = TDT^\top$  is a  $g$ -inverse of  $T^\top\Omega^{-1}T$ :

$$(A.4.5) \quad T^\top\Omega^{-1}TDT^\top\Omega^{-1}T \equiv T^\top\Omega^{-1}(\Omega - \Omega K^\top (K\Omega K^\top)^{-1} K\Omega)\Omega^{-1}T = T^\top\Omega^{-1}T.$$

To complete the proof of (5) we have to show that the expression  $T(T^\top\Omega^{-1}T)^{-1}T^\top$  does not depend on the choice of the  $g$ -inverse of  $T^\top\Omega^{-1}T$ . This follows from  $T(T^\top\Omega^{-1}T)^{-1}T^\top = \Omega\Omega^{-1}T(T^\top\Omega^{-1}T)^{-1}T^\top\Omega^{-1}\Omega$  and theorem A.5.10.

**THEOREM A.4.6.** *Given two matrices  $T$  and  $U$ . Then  $T \perp U$  if and only if for any  $D$  the following two statements are equivalent:*

$$(A.4.6) \quad TD = O$$

and

$$(A.4.7) \quad \text{For all } C \text{ which satisfy } CU = O \text{ follows } CD = O.$$

### A.5. Nonnegative Definite Symmetric Matrices

By definition, a symmetric matrix  $\Omega$  is nonnegative definite if  $\mathbf{a}^\top\Omega\mathbf{a} \geq 0$  for all vectors  $\mathbf{a}$ . It is positive definite if  $\mathbf{a}^\top\Omega\mathbf{a} > 0$  for all vectors  $\mathbf{a} \neq \mathbf{o}$ .

**THEOREM A.5.1.**  *$\Omega$  nonnegative definite symmetric if and only if it can be written in the form  $\Omega = A^\top A$  for some  $A$ .*

**THEOREM A.5.2.** *If  $\Omega$  is nonnegative definite, and  $\mathbf{a}^\top\Omega\mathbf{a} = 0$ , then already  $\Omega\mathbf{a} = \mathbf{o}$ .*

**THEOREM A.5.3.**  *$A$  is positive definite if and only if it is nonnegative definite and nonsingular.*

**THEOREM A.5.4.** *If the symmetric matrix  $A$  has a nnd  $g$ -inverse then  $A$  itself is also nnd.*

**THEOREM A.5.5.** *If  $\Omega$  and  $\Sigma$  are positive definite, then  $\Omega - \Sigma$  is positive (nonnegative) definite if and only if  $\Sigma^{-1} - \Omega^{-1}$  is.*

**THEOREM A.5.6.** *If  $\Omega$  and  $\Sigma$  are nonnegative definite, then  $\text{tr}(\Omega\Sigma) \geq 0$ .*

**PROBLEM 376.** *Prove theorem A.5.6.*

**ANSWER.** Find any factorization  $\Sigma = PP^\top$ . Then  $\text{tr}(\Omega\Sigma) = \text{tr}(P^\top\Omega P) \geq 0$ .

**THEOREM A.5.7.** *If  $\Omega$  is nonnegative definite symmetric, then*

$$(A.5.1) \quad (g^\top\Omega\mathbf{a})^2 \leq g^\top\Omega\mathbf{g} \mathbf{a}^\top\Omega\mathbf{a},$$

for arbitrary vectors  $\mathbf{a}$  and  $\mathbf{g}$ . Equality holds if and only if  $\Omega\mathbf{g}$  and  $\Omega\mathbf{a}$  are linearly dependent, i.e.,  $\alpha$  and  $\beta$  exist, not both zero, such that  $\Omega\mathbf{g}\alpha + \Omega\mathbf{a}\beta = \mathbf{o}$ .

Proof: First we will show that the condition for equality is sufficient. Therefore assume  $\Omega\mathbf{g}\alpha + \Omega\mathbf{a}\beta = \mathbf{o}$  for a certain  $\alpha$  and  $\beta$ , which are not both zero. Without loss of generality we can assume  $\alpha \neq 0$ . Then we can solve  $\mathbf{a}^\top\Omega\mathbf{g}\alpha + \mathbf{a}^\top\Omega\mathbf{a}\beta = 0$  to get  $\mathbf{a}^\top\Omega\mathbf{g} = -(\beta/\alpha)\mathbf{a}^\top\Omega\mathbf{a}$ , therefore the lefthand side of (A.5.1) is  $(\beta/\alpha)^2(\mathbf{a}^\top\Omega\mathbf{a})^2$ . Furthermore we can solve  $\mathbf{g}^\top\Omega\mathbf{g}\alpha + \mathbf{g}^\top\Omega\mathbf{a}\beta = 0$  to get  $\mathbf{g}^\top\Omega\mathbf{g} = -(\beta/\alpha)\mathbf{g}^\top\Omega\mathbf{a}$ , therefore the righthand side of (A.5.1) is  $(\beta/\alpha)^2(\mathbf{a}^\top\Omega\mathbf{a})^2$  as well—i.e., (A.5.1) holds with equality.

Secondly we will show that (A.5.1) holds in the general case and that, if it holds with equality,  $\Omega\mathbf{g}$  and  $\Omega\mathbf{a}$  are linearly dependent. We will split this second half of the proof into two substeps. First verify that (A.5.1) holds if  $\mathbf{g}^\top\Omega\mathbf{g} = 0$ . If this is the case, then already  $\Omega\mathbf{g} = \mathbf{o}$ , therefore the  $\Omega\mathbf{g}$  and  $\Omega\mathbf{a}$  are linearly dependent and by the first part of the proof, (A.5.1) holds with equality.

The second substep is the main part of the proof. Assume  $\mathbf{g}^\top\Omega\mathbf{g} \neq 0$ . Since  $\Omega$  is nonnegative definite, it follows

$$(A.5.2) \quad 0 \leq \left( \mathbf{a} - \mathbf{g} \frac{\mathbf{g}^\top\Omega\mathbf{a}}{\mathbf{g}^\top\Omega\mathbf{g}} \right)^\top \Omega \left( \mathbf{a} - \mathbf{g} \frac{\mathbf{g}^\top\Omega\mathbf{a}}{\mathbf{g}^\top\Omega\mathbf{g}} \right) = \mathbf{a}^\top\Omega\mathbf{a} - 2 \frac{(\mathbf{g}^\top\Omega\mathbf{a})^2}{\mathbf{g}^\top\Omega\mathbf{g}} + \frac{(\mathbf{g}^\top\Omega\mathbf{a})^2}{\mathbf{g}^\top\Omega\mathbf{g}} = \mathbf{a}^\top\Omega\mathbf{a} - \frac{(\mathbf{g}^\top\Omega\mathbf{a})^2}{\mathbf{g}^\top\Omega\mathbf{g}}$$

From this follows (A.5.1). If (A.5.2) is an equality, then already  $\Omega \left( \mathbf{a} - \mathbf{g} \frac{\mathbf{g}^\top\Omega\mathbf{a}}{\mathbf{g}^\top\Omega\mathbf{g}} \right) = \mathbf{o}$ , which means that  $\Omega\mathbf{g}$  and  $\Omega\mathbf{a}$  are linearly dependent.

**THEOREM A.5.8.** *In the situation of theorem A.5.7, one can take  $g$ -inverses without disturbing the inequality*

$$(A.5.3) \quad \frac{(\mathbf{g}^\top\Omega\mathbf{a})^2}{\mathbf{g}^\top\Omega\mathbf{g}} \leq \mathbf{a}^\top\Omega\mathbf{a}.$$

Equality holds if and only if a  $\gamma \neq 0$  exists with  $\Omega\mathbf{g} = \Omega\mathbf{a}\gamma$ .

PROBLEM 377. Show that if  $\Omega$  is nonnegative definite, then its elements satisfy

$$(A.5.4) \quad \omega_{ij}^2 \leq \omega_{ii}\omega_{jj}$$

ANSWER. Let  $\mathbf{a}$  and  $\mathbf{b}$  be the  $i$ th and  $j$ th unit vector. Then

$$(A.5.5) \quad \frac{(\mathbf{b}^\top \Omega \mathbf{a})^2}{\mathbf{b}^\top \Omega \mathbf{b}} \leq \max_g \frac{(g^\top \Omega \mathbf{a})^2}{g^\top \Omega g} = \mathbf{a}^\top \Omega \mathbf{a}.$$

□

PROBLEM 378. Assume  $\Omega$  nonnegative definite symmetric. If  $\mathbf{x}$  satisfies  $\mathbf{x} = \Omega \mathbf{a}$  for some  $\mathbf{a}$ , show that

$$(A.5.6) \quad \max_g \frac{(g^\top \mathbf{x})^2}{g^\top \Omega g} = \mathbf{x}^\top \Omega^- \mathbf{x}.$$

Furthermore show that equality holds if and only if  $\Omega \mathbf{g} = \mathbf{x} \gamma$  for some  $\gamma \neq 0$ .

ANSWER. From  $\mathbf{x} = \Omega \mathbf{a}$  follows  $g^\top \mathbf{x} = g^\top \Omega \mathbf{a}$  and  $\mathbf{x}^\top \Omega^- \mathbf{x} = \mathbf{a}^\top \Omega \mathbf{a}$ ; therefore it follows from theorem A.5.8.

□

PROBLEM 379. Assume  $\Omega$  nonnegative definite symmetric,  $\mathbf{x}$  satisfies  $\mathbf{x} = \Omega \mathbf{a}$  for some  $\mathbf{a}$ , and  $\mathbf{R}$  is such that  $\mathbf{R} \mathbf{x}$  is defined. Show that

$$(A.5.7) \quad \mathbf{x}^\top \mathbf{R}^\top (\mathbf{R} \Omega \mathbf{R}^\top)^- \mathbf{R} \mathbf{x} \leq \mathbf{x}^\top \Omega^- \mathbf{x}$$

ANSWER. Follows from

$$(A.5.8) \quad \max_h \frac{(h^\top \mathbf{R} \mathbf{x})^2}{h^\top \mathbf{R} \Omega \mathbf{R}^\top h} \leq \max_g \frac{(g^\top \mathbf{x})^2}{g^\top \Omega g}$$

because on the term on the lhs maximization is done over the smaller set of  $\mathbf{g}$  which have the form  $\mathbf{R} \mathbf{h}$ . An alternative proof would be to show that  $\Omega - \Omega \mathbf{r}^\top (\mathbf{R} \Omega \mathbf{R}^\top)^- \mathbf{R} \Omega$  is nnd (it has  $\Omega^-$  as  $\mathbf{g}$ -inverse).

□

PROBLEM 380. Assume  $\Omega$  nonnegative definite symmetric. Show that

$$(A.5.9) \quad \max_{\substack{g: \\ g = \Omega \mathbf{a} \\ \text{for some } \mathbf{a}}} \frac{(g^\top \mathbf{x})^2}{g^\top \Omega^- g} = \mathbf{x}^\top \Omega \mathbf{x}.$$

ANSWER. Since  $\mathbf{g} = \Omega \mathbf{a}$  for some  $\mathbf{a}$ , maximize over  $\mathbf{a}$  instead of  $\mathbf{g}$ . This reduces it to theorem A.5.8:

$$(A.5.10) \quad \max_{g: g = \Omega \mathbf{a} \text{ for some } \mathbf{a}} \frac{(g^\top \mathbf{x})^2}{g^\top \Omega^- g} = \max_{\mathbf{a}} \frac{(\mathbf{a}^\top \Omega \mathbf{x})^2}{\mathbf{a}^\top \Omega \mathbf{a}} = \mathbf{x}^\top \Omega \mathbf{x}$$

□

THEOREM A.5.9. Let  $\Omega$  be symmetric and nonnegative definite, and  $\mathbf{x}$  an arbitrary vector. Then  $\Omega - \mathbf{x} \mathbf{x}^\top$  is nonnegative definite if and only if the following conditions hold:  $\mathbf{x}$  can be written in the form  $\mathbf{x} = \Omega \mathbf{a}$  for some  $\mathbf{a}$ , and  $\mathbf{x}^\top \Omega^- \mathbf{x} \leq 1$  for one (and therefore for all)  $\mathbf{g}$ -inverses  $\Omega^-$  of  $\Omega$ .

PROBLEM 381. Prove theorem A.5.9.

ANSWER. Assume  $\mathbf{x} = \Omega \mathbf{a}$  and  $\mathbf{x}^\top \Omega^- \mathbf{x} = \mathbf{a}^\top \Omega \mathbf{a} \leq 1$ ; then for any  $\mathbf{g}$ ,  $g^\top (\Omega - \mathbf{x} \mathbf{x}^\top) g = g^\top \Omega g - g^\top \Omega \mathbf{a} \mathbf{a}^\top \Omega g \geq \mathbf{a}^\top \Omega \mathbf{a} g^\top \Omega g - g^\top \Omega \mathbf{a} \mathbf{a}^\top \Omega g \geq 0$  by theorem A.5.7.

Conversely, assume  $\mathbf{x}$  cannot be written in the form  $\mathbf{x} = \Omega \mathbf{a}$  for some  $\mathbf{a}$ ; then a  $\mathbf{g}$  exists with  $g^\top \Omega = \mathbf{o}^\top$  but  $g^\top \mathbf{x} \neq \mathbf{o}$ . Then  $g^\top (\Omega - \mathbf{x} \mathbf{x}^\top) g < 0$ , therefore not nnd.

Finally assume  $\mathbf{x}^\top \Omega^- \mathbf{x} = \mathbf{a}^\top \Omega \mathbf{a} > 1$ ; then  $\mathbf{a}^\top (\Omega - \mathbf{x} \mathbf{x}^\top) \mathbf{a} = \mathbf{a}^\top \Omega \mathbf{a} - (\mathbf{a}^\top \Omega \mathbf{a})^2 < 0$ , therefore again not nnd.

THEOREM A.5.10. If  $\Omega$  and  $\Sigma$  are nonnegative definite symmetric, and  $\mathbf{K}$  a matrix so that  $\Sigma \mathbf{K} \Omega$  is defined, then

$$(A.5.11) \quad \mathbf{K} \Omega = (\mathbf{K} \Omega \mathbf{K}^\top + \Sigma) (\mathbf{K} \Omega \mathbf{K}^\top + \Sigma)^- \mathbf{K} \Omega.$$

Furthermore,  $\Omega \mathbf{K}^\top (\mathbf{K} \Omega \mathbf{K}^\top + \Sigma)^- \mathbf{K} \Omega$  is independent of the choice of  $\mathbf{g}$ -inverses.

PROBLEM 382. Prove theorem A.5.10.

ANSWER. To see that (A.5.11) is a special case of (A.3.3), take any  $\mathbf{Q}$  with  $\Omega = \mathbf{Q} \mathbf{Q}^\top$  and with  $\Sigma = \mathbf{P} \mathbf{P}^\top$  and define  $\mathbf{A} = \begin{bmatrix} \mathbf{K} \mathbf{Q} & \mathbf{P} \end{bmatrix}$ . The independence of the choice of  $\mathbf{g}$ -inverses follows from theorem A.3.1 together with (A.5.11).

The following was apparently first shown in [Alb69] for the special case of the Moore-Penrose pseudoinverse:

THEOREM A.5.11. The symmetric partitioned matrix  $\Omega = \begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{yz} & \Omega_{zz} \end{bmatrix}$  is nonnegative definite if and only if the following conditions hold:

$$(A.5.12) \quad \Omega_{yy} \text{ and } \Omega_{zz.y} := \Omega_{zz} - \Omega_{yz}^\top \Omega_{yy}^- \Omega_{yz} \text{ are both nonnegative definite, and}$$

$$(A.5.13) \quad \Omega_{yz} = \Omega_{yy} \Omega_{yy}^- \Omega_{yz}$$

Reminder: It follows from theorem A.3.1 that (A.5.13) holds for some  $\mathbf{g}$ -inverse if and only if it holds for all, and that, if it holds,  $\Omega_{zz.y}$  is independent of the choice of the  $\mathbf{g}$ -inverse.

Proof of theorem A.5.11: First we prove the necessity of the three conditions in the theorem. If the symmetric partitioned matrix  $\Omega$  is nonnegative definite

there exists a  $\mathbf{R}$  with  $\Omega = \mathbf{R}^\top \mathbf{R}$ . Write  $\mathbf{R} = \begin{bmatrix} \mathbf{R}_y & \mathbf{R}_z \end{bmatrix}$  to get  $\begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{yz} & \Omega_{zz} \end{bmatrix} = \begin{bmatrix} \mathbf{R}_y^\top & \mathbf{R}_z^\top \\ \mathbf{R}_y & \mathbf{R}_z \end{bmatrix} \begin{bmatrix} \mathbf{R}_y & \mathbf{R}_z \\ \mathbf{R}_y & \mathbf{R}_z \end{bmatrix}$ .

$\Omega_{yy}$  is nonnegative definite because it is equal to  $\mathbf{R}_y^\top \mathbf{R}_y$ , and



(A.5.13) follows from (A.5.11):  $\Omega_{yy}\Omega_{yy}^{-1}\Omega_{yz} = R_y^T R_y (R_y^T R_y)^{-1} R_y^T R_z = R_y^T R_z = \Omega_{yz}$ . To show that  $\Omega_{zz.y}$  is nonnegative definite, define  $S = (I - R_y (R_y^T R_y)^{-1} R_y^T) R_z$ . Then  $S^T S = R_z^T (I - R_y (R_y^T R_y)^{-1} R_y^T) R_z = \Omega_{zz.y}$ .

To show sufficiency of the three conditions of theorem A.5.11, assume the symmetric  $\begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{yz}^T & \Omega_{zz} \end{bmatrix}$  satisfies them. Pick two matrices  $Q$  and  $S$  so that  $\Omega_{yy} = Q^T Q$  and  $\Omega_{zz.y} = S^T S$ . Then

$$\begin{bmatrix} \Omega_{yy} & \Omega_{yz} \\ \Omega_{yz}^T & \Omega_{zz} \end{bmatrix} = \begin{bmatrix} Q^T & O \\ \Omega_{yz}^T \Omega_{yy}^{-1} Q^T & S^T \end{bmatrix} \begin{bmatrix} Q & Q \Omega_{yy}^{-1} \Omega_{yz} \\ O & S \end{bmatrix},$$

therefore nonnegative definite.

PROBLEM 383. [SM86, A 3.2/11] Given a positive definite matrix  $Q$  and a positive definite  $\tilde{Q}$  with  $Q^* = Q - \tilde{Q}$  nonnegative definite.

- a. Show that  $\tilde{Q} - \tilde{Q}Q^{-1}\tilde{Q}$  is nonnegative definite.

ANSWER. We know that  $\tilde{Q}^{-1} - Q^{*-1}$  is nnd, therefore  $\tilde{Q}\tilde{Q}^{-1}\tilde{Q} - \tilde{Q}Q^{*-1}\tilde{Q}$  nnd.  $\square$

- b. This part is more difficult: Show that also  $Q^* - Q^*Q^{-1}Q^*$  is nonnegative definite.

ANSWER. We will write it in a symmetric form from which it is obvious that it is nonnegative definite:

$$(A.5.14) \quad Q^* - Q^*Q^{-1}Q^* = Q^* - Q^*(\tilde{Q} + Q^*)^{-1}Q^*$$

$$(A.5.15) \quad = Q^*(\tilde{Q} + Q^*)^{-1}(\tilde{Q} + Q^* - Q^*) = Q^*(\tilde{Q} + Q^*)^{-1}\tilde{Q}$$

$$(A.5.16) \quad = \tilde{Q}(\tilde{Q} + Q^*)^{-1}(\tilde{Q} + Q^*)\tilde{Q}^{-1}Q^*(\tilde{Q} + Q^*)^{-1}\tilde{Q}$$

$$(A.5.17) \quad = \tilde{Q}Q^{-1}(Q^* + Q^*\tilde{Q}^{-1}Q^*)Q^{-1}\tilde{Q}.$$

$\square$

PROBLEM 384. Given the vector  $h \neq o$ . For which values of the scalar  $\gamma$  is the matrix  $I - \frac{hh^T}{\gamma}$  singular, nonsingular, nonnegative definite, a projection matrix, orthogonal?

ANSWER. It is nnd iff  $\gamma \geq h^T h$ , because of theorem A.5.9. One easily verifies that it is orthogonal iff  $\gamma = h^T h/2$ , and it is a projection matrix iff  $\gamma = h^T h$ . Now let us prove that it is singular iff  $\gamma = h^T h$ : if this condition holds, then the matrix annuls  $h$ ; now assume the condition does not hold, i.e.,  $\gamma \neq h^T h$ , and take any  $x$  with  $(I - \frac{hh^T}{\gamma})x = o$ . It follows  $x = h\alpha$  where  $\alpha = h^T x/\gamma$ , therefore  $(I - \frac{hh^T}{\gamma})x = h\alpha(1 - h^T h/\gamma)$ . Since  $h \neq o$  and  $1 - h^T h/\gamma \neq 0$  this can only be the null vector if  $\alpha = 0$ .  $\square$

### A.6. Projection Matrices

PROBLEM 385. Show that  $X(X^T X)^{-1}X^T$  is the projection matrix on the range space  $R[X]$  of  $X$ , i.e., on the space spanned by the columns of  $X$ . This is true whether or not  $X$  has full column rank.

ANSWER. Idempotence requires theorem A.3.2, and symmetry the invariance under choice of g-inverse. Furthermore one has to show  $X(X^T X)^{-1}Xa = a$  holds if and only if  $a = Xb$  for some  $b$ .  $\Rightarrow$  is clear, and  $\Leftarrow$  follows from theorem A.3.2.

THEOREM A.6.1. Let  $P$  and  $Q$  be projection matrices, i.e., both are symmetric and idempotent. Then the following five conditions are equivalent, each meaning that the space on which  $P$  projects is a subspace of the space on which  $Q$  projects:

$$(A.6.1) \quad R[P] \subset R[Q]$$

$$(A.6.2) \quad QP = P$$

$$(A.6.3) \quad PQ = P$$

$$(A.6.4) \quad Q - P \text{ projection matrix}$$

$$(A.6.5) \quad Q - P \text{ nonnegative definite.}$$

(A.6.2) is geometrically trivial. It means: if one first projects on a certain space and then on a larger space which contains the first space as a subspace, then nothing happens under this second projection because one is already in the larger space. (A.6.3) is geometrically not trivial and worth remembering: if one first projects on a certain space, and then on a smaller space which is a subspace of the first space, then the result is the same as if one had projected directly on the smaller space. (A.6.4) means: the difference  $Q - P$  is the projection on the orthogonal complement of  $R[P]$  in  $R[Q]$ . And (A.6.5) means: the projection of a vector on the smaller space cannot be longer than that on the larger space.

PROBLEM 386. Prove theorem A.6.1.

ANSWER. Instead of going in a circle it is more natural to show (A.6.1)  $\iff$  (A.6.2)  $\iff$  (A.6.3)  $\iff$  (A.6.2) and then go in a circle for the remaining conditions: (A.6.2), (A.6.3)  $\implies$  (A.6.4)  $\implies$  (A.6.5).

(A.6.1)  $\implies$  (A.6.2):  $R[P] \subset R[Q]$  means that for every  $c$  exists a  $d$  with  $Pc = Qd$ . Therefore for all  $c$  follows  $QPc = QQd = Qd = Pc$ , i.e.,  $QP = P$ .

(A.6.2)  $\implies$  (A.6.1): if  $Pc = QPc$  for all  $c$ , then clearly  $R[P] \subset R[Q]$ .

(A.6.2)  $\implies$  (A.6.3) by symmetry of  $P$  and  $Q$ : If  $QP = P$  then  $PQ = P^T Q^T = (QP)^T = P^T = P$ .

(A.6.3)  $\implies$  (A.6.2) follows in exactly the same way: If  $PQ = P$  then  $QP = Q^T P^T = (PQ)^T = P^T = P$ .

(A.6.2), (A.6.3)  $\implies$  (A.6.4): Symmetry of  $Q - P$  clear, and  $(Q - P)(Q - P) = Q - P - P + P^2 = Q - P$ .

$$(A.6.4) \Rightarrow (A.6.5): \mathbf{c}^\top(\mathbf{Q} - \mathbf{P})\mathbf{c} = \mathbf{c}^\top(\mathbf{Q} - \mathbf{P})^\top(\mathbf{Q} - \mathbf{P})\mathbf{c} \geq 0.$$

(A.6.5)  $\Rightarrow$  (A.6.3): First show that, if  $\mathbf{Q} - \mathbf{P}$  nnd, then  $\mathbf{Q}\mathbf{c} = \mathbf{o}$  implies  $\mathbf{P}\mathbf{c} = \mathbf{o}$ . Proof: from  $\mathbf{Q} - \mathbf{P}$  nnd and  $\mathbf{Q}\mathbf{c} = \mathbf{o}$  follows  $0 \leq \mathbf{c}^\top(\mathbf{Q} - \mathbf{P})\mathbf{c} = -\mathbf{c}^\top\mathbf{P}\mathbf{c} \leq 0$ , therefore equality throughout, i.e.,  $0 = \mathbf{c}^\top\mathbf{P}\mathbf{c} = \mathbf{c}^\top\mathbf{P}^\top\mathbf{P}\mathbf{c} = \|\mathbf{P}\mathbf{c}\|^2$  and therefore  $\mathbf{P}\mathbf{c} = \mathbf{o}$ . Secondly: this is also true for matrices:  $\mathbf{Q}\mathbf{C} = \mathbf{O}$  implies  $\mathbf{P}\mathbf{C} = \mathbf{O}$ , since it is valid for every column of  $\mathbf{C}$ . Thirdly: Since  $\mathbf{Q}(\mathbf{I} - \mathbf{Q}) = \mathbf{O}$ , it follows  $\mathbf{P}(\mathbf{I} - \mathbf{Q}) = \mathbf{O}$ , which is (A.6.3).  $\square$

PROBLEM 387. If  $\mathbf{Y} = \mathbf{X}\mathbf{A}$  for some  $\mathbf{A}$ , show that  $\mathbf{Y}(\mathbf{Y}^\top\mathbf{Y})^{-1}\mathbf{Y}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top = \mathbf{Y}(\mathbf{Y}^\top\mathbf{Y})^{-1}\mathbf{Y}^\top$ .

ANSWER.  $\mathbf{Y} = \mathbf{X}\mathbf{A}$  means that every column of  $\mathbf{Y}$  is a linear combination of columns of  $\mathbf{A}$ :

$$(A.6.6) \quad \begin{bmatrix} \mathbf{y}_1 & \cdots & \mathbf{y}_m \end{bmatrix} = \mathbf{X} \begin{bmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_m \end{bmatrix} = \begin{bmatrix} \mathbf{X}\mathbf{a}_1 & \cdots & \mathbf{X}\mathbf{a}_m \end{bmatrix}.$$

Therefore geometrically the statement follows from the fact shown in Problem 385 that the above matrices are projection matrices on the columnn spaces. But it can also be shown algebraically:  $\mathbf{Y}(\mathbf{Y}^\top\mathbf{Y})^{-1}\mathbf{Y}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top = \mathbf{Y}(\mathbf{Y}^\top\mathbf{Y})^{-1}\mathbf{A}^\top\mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top = \mathbf{Y}(\mathbf{Y}^\top\mathbf{Y})^{-1}\mathbf{Y}^\top$ .  $\square$

PROBLEM 388. (Not eligible for in-class exams) Let  $\mathbf{Q}$  be a projection matrix (i.e., a symmetric and idempotent matrix) with the property that  $\mathbf{Q} = \mathbf{X}\mathbf{A}\mathbf{X}^\top$  for some  $\mathbf{A}$ . Define  $\tilde{\mathbf{X}} = (\mathbf{I} - \mathbf{Q})\mathbf{X}$ . Then

$$(A.6.7) \quad \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top = \tilde{\mathbf{X}}(\tilde{\mathbf{X}}^\top\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}^\top + \mathbf{Q}.$$

Hint: this can be done through a geometric argument. If you want to do it algebraically, you might want to use the fact that  $(\mathbf{X}^\top\mathbf{X})^{-1}$  is also a g-inverse of  $\tilde{\mathbf{X}}^\top\tilde{\mathbf{X}}$ .

ANSWER. Geometric argument:  $\mathbf{Q}$  is a projector on a subspace of the range space of  $\mathbf{X}$ . The columns of  $\tilde{\mathbf{X}}$  are projections of the columns of  $\mathbf{X}$  on the orthogonal complement of the space on which  $\mathbf{Q}$  projects. The equation which we have to prove shows therefore that the projection on the column space of  $\mathbf{X}$  is the sum of the projections on the space  $\mathbf{Q}$  projects on plus the projection on the orthogonal complement of that space in  $\mathbf{X}$ .

Now an algebraic proof: First let us show that  $(\mathbf{X}^\top\mathbf{X})^{-1}$  is a g-inverse of  $\tilde{\mathbf{X}}^\top\tilde{\mathbf{X}}$ , i.e., let us evaluate

$$(A.6.8) \quad \mathbf{X}^\top(\mathbf{I} - \mathbf{Q})\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top(\mathbf{I} - \mathbf{Q})\mathbf{X} = \mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X} - \mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{Q}\mathbf{X} - \mathbf{X}^\top\mathbf{Q}\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}$$

$$(A.6.9) \quad = \mathbf{X}^\top\mathbf{X} - \mathbf{X}^\top\mathbf{Q}\mathbf{X} - \mathbf{X}^\top\mathbf{Q}\mathbf{X} + \mathbf{X}^\top\mathbf{Q}\mathbf{X} = \mathbf{X}^\top(\mathbf{I} - \mathbf{Q})\mathbf{X}.$$

Only for the fourth term did we need the condition  $\mathbf{Q} = \mathbf{X}\mathbf{A}\mathbf{X}^\top$ :

$$(A.6.10) \quad \mathbf{X}^\top\mathbf{X}\mathbf{A}\mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}\mathbf{A}\mathbf{X}^\top\mathbf{X} = \mathbf{X}^\top\mathbf{X}\mathbf{A}\mathbf{X}^\top\mathbf{X}\mathbf{A}\mathbf{X}^\top\mathbf{X} = \mathbf{X}^\top\mathbf{Q}\mathbf{Q}\mathbf{X} = \mathbf{X}^\top\mathbf{X}.$$

Using this g-inverse we have

$$(A.6.11) \quad \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top - \tilde{\mathbf{X}}(\tilde{\mathbf{X}}^\top\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}^\top = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top - (\mathbf{I} - \mathbf{Q})\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top(\mathbf{I} - \mathbf{Q}) =$$

$$(A.6.12) \quad = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top - \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{Q} + \mathbf{Q}\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top - \mathbf{Q}\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{Q} = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top \quad \square$$

PROBLEM 389. Given any projection matrix  $\mathbf{P}$ . Show that its  $i$ th diagonal element can be written

$$(A.6.13) \quad p_{ii} = \sum_j p_{ij}^2.$$

ANSWER. From idempotence  $\mathbf{P} = \mathbf{P}\mathbf{P}$  follows  $p_{ii} = \sum_j p_{ij}p_{ji}$ , now use symmetry to (A.6.13).

### A.7. Determinants

THEOREM A.7.1. The determinant of a block-triangular matrix is the product of the determinants of the blocks in the diagonal. In other words,

$$(A.7.1) \quad \begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{O} & \mathbf{D} \end{vmatrix} = |\mathbf{A}| |\mathbf{D}|$$

For the proof recall the definition of a determinant. A mapping  $\pi: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  is called a permutation if and only if it is one-to-one if and only if it is onto. Permutations can be classified as even or odd according to whether they can be written as the product of an even or odd number of transpositions. Then the determinant is defined as

$$(A.7.2) \quad \det(\mathbf{A}) = \sum_{\pi: \pi \text{ even}} a_{1\pi(1)} \cdots a_{n\pi(n)} - \sum_{\pi: \pi \text{ odd}} a_{1\pi(1)} \cdots a_{n\pi(n)}$$

Now assume  $\mathbf{A}$  is  $m \times m$ ,  $1 \leq m < n$ . If a  $j \leq m$  exists with  $\pi(j) > m$  then not all  $i \leq m$  can be images of other points  $j \leq m$ , i.e., there must be at least one  $j > m$  with  $\pi(j) \leq m$ . Therefore, in a block triangular matrix in which all  $a_{ij} = 0$  for  $i \leq m, j > m$ , only those permutations give a nonzero product which remain within the two submatrices straddling the diagonal.

THEOREM A.7.2. If  $\mathbf{B} = \mathbf{A}\mathbf{A}^{-1}\mathbf{B}$ , then the following identity is valid between determinants:

$$(A.7.3) \quad \begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{vmatrix} = |\mathbf{A}| |\mathbf{E}| \quad \text{where } \mathbf{E} = \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}.$$

Proof: Postmultiply by a matrix whose determinant, by lemma A.7.1, is one and then apply lemma A.7.1 once more:

$$(A.7.4) \quad \begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{vmatrix} = \begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{vmatrix} \begin{vmatrix} \mathbf{I} & -\mathbf{A}^{-1}\mathbf{B} \\ \mathbf{O} & \mathbf{I} \end{vmatrix} = \begin{vmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{C} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{vmatrix} = |\mathbf{A}| |\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}|.$$

PROBLEM 390. Show the following counterpart of theorem A.7.2: If  $C = DD^{-1}C$ , then the following identity is valid between determinants:

$$(A.7.5) \quad \begin{vmatrix} A & B \\ C & D \end{vmatrix} = |A - BD^{-1}C| |D|.$$

ANSWER.

$$(A.7.6) \quad \begin{vmatrix} A & B \\ C & D \end{vmatrix} = \begin{vmatrix} A & B \\ C & D \\ -D^{-1}C & I \end{vmatrix} = \begin{vmatrix} A - BD^{-1}C & B \\ O & D \end{vmatrix} = |A - BD^{-1}C| |D|.$$

□

PROBLEM 391. Show that whenever  $BC$  and  $CB$  are defined, it follows  $|I - BC| = |I - CB|$ .

ANSWER. Set  $A = I$  and  $D = I$  in (A.7.3) and (A.7.5). □

THEOREM A.7.3. Assume that  $d = WW^{-1}d$ . Then

$$(A.7.7) \quad \det(W + \alpha \cdot dd^T) = \det(W)(1 + \alpha d^T W^{-1}d).$$

Proof: If  $\alpha = 0$ , then there is nothing to prove. Otherwise look at the determinant of the matrix

$$(A.7.8) \quad H = \begin{bmatrix} W & d \\ d^T & -1/\alpha \end{bmatrix}$$

Equations (A.7.3) and (A.7.5) give two expressions for it:

$$(A.7.9) \quad \det(H) = \det(W)(-1/\alpha - d^T W^{-1}d) = -\frac{1}{\alpha} \det(W + \alpha dd^T).$$

### A.8. More About Inverses

PROBLEM 392. Given a partitioned matrix  $\begin{bmatrix} A & B \\ C & D \end{bmatrix}$  which satisfies  $B = AA^{-1}B$  and  $C = CA^{-1}A$ . (These conditions hold for instance, due to theorem A.5.11, if  $\begin{bmatrix} A & B \\ C & D \end{bmatrix}$  is nonnegative definite symmetric, but it also holds in the nonsymmetric case if  $A$  is nonsingular, which by theorem A.7.2 is the case if the whole partitioned matrix is nonsingular.) Define  $E = D - CA^{-1}B$ ,  $F = A^{-1}B$ , and  $G = CA^{-1}$ .

• a. Prove that in terms of  $A$ ,  $E$ ,  $F$ , and  $G$ , the original matrix can be written as

$$(A.8.1) \quad \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix}$$

(this is trivial), and that (this is the nontrivial part)

$$(A.8.2) \quad \begin{bmatrix} A^{-1} + FE^{-1}G & -FE^{-1} \\ -E^{-1}G & E^{-1} \end{bmatrix} \text{ is a } g\text{-inverse of } \begin{bmatrix} A & B \\ C & D \end{bmatrix}.$$

ANSWER. This here is not the shortest proof because I was still wondering if it could be formulated in a more general way. Multiply out but do not yet use the conditions  $B = AA^{-1}B$  and  $C = CA^{-1}A$ :

$$(A.8.3) \quad \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} A^{-1} + FE^{-1}G & -FE^{-1} \\ -E^{-1}G & E^{-1} \end{bmatrix} = \begin{bmatrix} AA^{-1} - (I - AA^{-1})BE^{-1}G & (I - AA^{-1})BE^{-1} \\ (I - EE^{-1})G & EE^{-1} \end{bmatrix}$$

and

$$(A.8.4) \quad \begin{bmatrix} AA^{-1} - (I - AA^{-1})BE^{-1}G & (I - AA^{-1})BE^{-1} \\ (I - EE^{-1})G & EE^{-1} \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} A + (I - AA^{-1})BE^{-1}C(I - A^{-1}A) & B - (I - AA^{-1})B(I - E^{-1}E) \\ C - (I - EE^{-1})C(I - A^{-1}A) & D \end{bmatrix}$$

One sees that not only the conditions  $B = AA^{-1}B$  and  $C = CA^{-1}A$ , but also the conditions  $B = AA^{-1}B$  and  $C = EE^{-1}C$ , or alternatively the conditions  $B = BE^{-1}E$  and  $C = CA^{-1}A$  imply the statement. I think one can also work with the conditions  $AA^{-1}B = BD^{-1}D$  and  $DD^{-1}C = CA^{-1}A$ . Note that the lower right partition is  $D$  no matter what.

• b. If  $\begin{bmatrix} U & V \\ W & X \end{bmatrix}$  is a  $g$ -inverse of  $\begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix}$ , show that  $X$  is a  $g$ -inverse of  $E$ .

ANSWER. The  $g$ -inverse condition means

$$(A.8.5) \quad \begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix} \begin{bmatrix} U & V \\ W & X \end{bmatrix} \begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix} = \begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix}$$

The first matrix product evaluated is

$$(A.8.6) \quad \begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix} \begin{bmatrix} U & V \\ W & X \end{bmatrix} = \begin{bmatrix} AU + AFW & AV + AFX \\ GAU + EW + GAFW & GAV + EX + GAFX \end{bmatrix}$$

The  $g$ -inverse condition means therefore

$$(A.8.7) \quad \begin{bmatrix} AU + AFW & AV + AFX \\ GAU + EW + GAFW & GAV + EX + GAFX \end{bmatrix} \begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix} = \begin{bmatrix} A & AF \\ GA & E + GAF \end{bmatrix}$$

For the upper left partition this means  $AUA + AFWA + AVGA + AFXGA = A$ , and for the upper right partition it means  $AUAF + AFWAF + AVE + AVGAF + AFXE + AFXGAF = A$ . Postmultiply the upper left equation by  $F$  and subtract from the upper right to get  $AVB + AFXE = O$ . For the lower left we get  $GAUA + EWA + GAFWA + GAVGA + EXGA + GAFXGA = GA$ . Premultiplication of the upper left equation by  $G$  and subtraction gives  $EW + EXGA = O$ . For the lower right corner we get  $GAUAF + EWAF + GAFWAF + GAVAF + EXE + GAFXE + GAVGAF + EXGAF + GAFXGAF = E + GAF$ . Since  $AVE + AFXE = E$  and  $EWA + EXGA = O$ , this simplifies to  $GAUAF + GAFWAF + EXE + GAVGAF = E + GAF$ .

$GAFXGAF = E + GAF$ . And if one premultiplies the upper right corner by  $G$  and postmultiplies it by  $F$  and subtracts it from this one gets  $EXE = E$ .  $\square$

PROBLEM 393. Show that a  $g$ -inverse of the matrix

$$(A.8.8) \quad \begin{bmatrix} \mathbf{X}_1^\top \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{X}_2 \\ \mathbf{X}_2^\top \mathbf{X}_1 & \mathbf{X}_2^\top \mathbf{X}_2 \end{bmatrix}$$

has the form

$$(A.8.9) \quad \begin{bmatrix} (\mathbf{X}_1^\top \mathbf{X}_1)^- + \mathbf{D}_1^\top \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^- \mathbf{X}_2^\top \mathbf{D}_1 & -\mathbf{D}_1^\top \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^- \\ -(\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^- \mathbf{X}_2^\top \mathbf{D}_1 & (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^- \end{bmatrix}$$

where  $\mathbf{M}_1 = \mathbf{I} - \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^- \mathbf{X}_1^\top$  and  $\mathbf{D}_1 = \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^-$ .

ANSWER. Either show it by multiplying it out, or apply Problem 392.  $\square$

PROBLEM 394. Show that the following are  $g$ -inverses:

$$(A.8.10) \quad \begin{bmatrix} \mathbf{I} & \mathbf{X} \\ \mathbf{X}^\top & \mathbf{X}^\top \mathbf{X} \end{bmatrix}^- = \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \quad \begin{bmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \\ \mathbf{X} & \mathbf{I} \end{bmatrix}^- = \begin{bmatrix} (\mathbf{X}^\top \mathbf{X})^- & \mathbf{O} \\ \mathbf{O} & \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X} \end{bmatrix}$$

ANSWER. Either do it by multiplying it out, or apply problem 392.  $\square$

PROBLEM 395. Assume again  $\mathbf{B} = \mathbf{A}\mathbf{A}^- \mathbf{B}$  and  $\mathbf{C} = \mathbf{C}\mathbf{A}\mathbf{A}^-$ , but assume this time that  $\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}$  nonsingular. Then  $\mathbf{A}$  is nonsingular,

$$(A.8.11) \quad \text{and if } \begin{bmatrix} \mathbf{P} & \mathbf{Q} \\ \mathbf{R} & \mathbf{S} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1}, \quad \text{then the determinant } \begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{vmatrix} = \frac{|\mathbf{A}|}{|\mathbf{S}|}.$$

ANSWER. The determinant is, by (A.7.3),  $|\mathbf{A}| |\mathbf{E}|$  where  $\mathbf{E} = \mathbf{D} - \mathbf{C}\mathbf{A}^- \mathbf{B}$ . By assumption, this determinant is nonzero, therefore also  $|\mathbf{A}|$  and  $|\mathbf{E}|$  are nonzero, i.e.,  $\mathbf{A}$  and  $\mathbf{E}$  are nonsingular. Therefore (A.8.2) reads

$$(A.8.12) \quad \begin{bmatrix} \mathbf{P} & \mathbf{Q} \\ \mathbf{R} & \mathbf{S} \end{bmatrix} = \begin{bmatrix} \mathbf{A}^{-1} + \mathbf{F}\mathbf{E}^{-1}\mathbf{G} & -\mathbf{F}\mathbf{E}^{-1} \\ -\mathbf{E}^{-1}\mathbf{G} & \mathbf{E}^{-1} \end{bmatrix},$$

i.e.,  $\mathbf{S} = \mathbf{E}^{-1} = (\mathbf{D} - \mathbf{C}\mathbf{A}^- \mathbf{B})^{-1}$ . hence  $|\mathbf{A}| |\mathbf{E}| = |\mathbf{A}| / |\mathbf{S}|$ .  $\square$

THEOREM A.8.1. Given a  $m \times n$  matrix  $\mathbf{A}$ , a  $m \times h$  matrix  $\mathbf{B}$ , a  $k \times n$  matrix  $\mathbf{C}$ , and a  $k \times h$  matrix  $\mathbf{D}$  satisfying  $\mathbf{A}\mathbf{A}^- \mathbf{B} = \mathbf{B}\mathbf{D}^- \mathbf{D}$  and  $\mathbf{D}\mathbf{D}^- \mathbf{C} = \mathbf{C}\mathbf{A}^- \mathbf{A}$ . Then the following are  $g$ -inverses:

$$(A.8.13) \quad (\mathbf{A} + \mathbf{B}\mathbf{D}^- \mathbf{C})^- = \mathbf{A}^- - \mathbf{A}^- \mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^- \mathbf{B})^- \mathbf{C}\mathbf{A}^-$$

$$(A.8.14) \quad (\mathbf{D} + \mathbf{C}\mathbf{A}^- \mathbf{B})^- = \mathbf{D}^- - \mathbf{D}^- \mathbf{C}(\mathbf{A} + \mathbf{B}\mathbf{D}^- \mathbf{C})^- \mathbf{B}\mathbf{D}^-.$$

PROBLEM 396. Prove theorem A.8.1.

ANSWER. Proof: Define  $\mathbf{E} = \mathbf{D} + \mathbf{C}\mathbf{A}^- \mathbf{B}$ . Then it follows from the assumptions that

$$(A.8.15) \quad (\mathbf{A} + \mathbf{B}\mathbf{D}^- \mathbf{C})(\mathbf{A}^- - \mathbf{A}^- \mathbf{B}\mathbf{E}^- \mathbf{C}\mathbf{A}^-) = \mathbf{A}\mathbf{A}^- - \mathbf{B}\mathbf{D}^- \mathbf{D}\mathbf{E}^- \mathbf{C}\mathbf{A}^- + \mathbf{B}\mathbf{D}^- \mathbf{C}\mathbf{A}^- - \mathbf{B}\mathbf{D}^- \mathbf{C}\mathbf{A}^- \mathbf{B}\mathbf{E}^- \mathbf{C}\mathbf{A}^-$$

$$(A.8.16) \quad = \mathbf{A}\mathbf{A}^- + \mathbf{B}\mathbf{D}^- (\mathbf{I} - \mathbf{E}\mathbf{E}^-) \mathbf{C}\mathbf{A}^-$$

Since  $\mathbf{A}\mathbf{A}^- (\mathbf{A} + \mathbf{B}\mathbf{D}^- \mathbf{C}) = \mathbf{A} + \mathbf{B}\mathbf{D}^- \mathbf{C}$ , we have to show that the second term on the rhs. ann

$$(A.8.17) \quad \mathbf{B}\mathbf{D}^- (\mathbf{I} - \mathbf{E}\mathbf{E}^-) \mathbf{C}\mathbf{A}^- (\mathbf{A} + \mathbf{B}\mathbf{D}^- \mathbf{C}) =$$

$$(A.8.18) \quad = \mathbf{B}\mathbf{D}^- \mathbf{C}\mathbf{A}^- \mathbf{A} + \mathbf{B}\mathbf{D}^- \mathbf{C}\mathbf{A}^- \mathbf{B}\mathbf{D}^- \mathbf{C} - \mathbf{B}\mathbf{D}^- \mathbf{E}\mathbf{E}^- \mathbf{C}\mathbf{A}^- \mathbf{A} - \mathbf{B}\mathbf{D}^- \mathbf{E}\mathbf{E}^- \mathbf{C}\mathbf{A}^- \mathbf{B}\mathbf{D}^- \mathbf{C}$$

$$(A.8.19) \quad = \mathbf{B}\mathbf{D}^- (\mathbf{D} + \mathbf{C}\mathbf{A}^- \mathbf{B} - \mathbf{E}\mathbf{E}^- \mathbf{D} - \mathbf{E}\mathbf{E}^- \mathbf{C}\mathbf{A}^- \mathbf{B}) \mathbf{D}^- \mathbf{C} = \mathbf{B}\mathbf{D}^- (\mathbf{E} - \mathbf{E}\mathbf{E}^- \mathbf{E}) \mathbf{D}^- \mathbf{C} = \mathbf{O}.$$

THEOREM A.8.2. (Sherman-Morrison-Woodbury theorem) Given a  $m \times n$  matrix  $\mathbf{A}$ , a  $m \times 1$  vector  $\mathbf{b}$  satisfying  $\mathbf{A}\mathbf{A}^- \mathbf{b} = \mathbf{b}$ , a  $n \times 1$  vector  $\mathbf{c}$  satisfying  $\mathbf{c}^\top \mathbf{A}\mathbf{A}^- = \mathbf{c}^\top$  and a scalar  $\delta$ . If  $\mathbf{A}^-$  is a  $g$ -inverse of  $\mathbf{A}$ , then

$$(A.8.20) \quad \mathbf{A}^- - \frac{\mathbf{A}^- \mathbf{b} \mathbf{c}^\top \mathbf{A}^-}{\mathbf{c}^\top \mathbf{A}^- \mathbf{b} + \delta} \quad \text{is a } g\text{-inverse of } \mathbf{A} + \frac{\mathbf{b} \mathbf{c}^\top}{\delta}$$

PROBLEM 397. Prove theorem A.8.2.

ANSWER. It is a special case of theorem A.8.1.

THEOREM A.8.3. For any symmetric nonnegative definite  $r \times r$  matrix  $\mathbf{A}$ ,

$$(A.8.21) \quad (\det \mathbf{A}) e^{-(\text{tr } \mathbf{A})} \leq e^{-r},$$

with equality holding if and only if  $\mathbf{A} = \mathbf{I}$ .

PROBLEM 398. Prove Theorem A.8.3. Hint: Let  $\lambda_1, \dots, \lambda_r$  be the eigenvalues of  $\mathbf{A}$ . Then  $\det \mathbf{A} = \prod_i \lambda_i$ , and  $\text{tr } \mathbf{A} = \sum_i \lambda_i$ .

ANSWER. Therefore the inequality reads

$$(A.8.22) \quad \prod_{i=1}^r \lambda_i e^{-\lambda_i} \leq e^{-r}$$

For this it is sufficient to show for each value of  $\lambda$

$$(A.8.23) \quad \lambda e^{-\lambda} \leq e^{-1},$$

which follows immediately by taking the derivatives:  $e^{-\lambda} - \lambda e^{-\lambda} = 0$  gives  $\lambda = 1$ . The matrix with all eigenvalues being equal to 1 is the identity matrix.

### A.9. Eigenvalues and Singular Value Decomposition

Every symmetric matrix  $\mathbf{B}$  has real eigenvalues and a system of orthogonal eigenvectors which span the whole space. If one normalizes these eigenvectors and combines them as row vectors into a matrix  $\mathbf{T}$ , then orthonormality means  $\mathbf{TT}^\top = \mathbf{I}$ , and since  $\mathbf{T}$  is square,  $\mathbf{TT}^\top = \mathbf{I}$  also implies  $\mathbf{T}^\top\mathbf{T} = \mathbf{I}$ , i.e.,  $\mathbf{T}$  is an orthogonal matrix. The existence of a complete set of real eigenvectors is therefore equivalent to the following matrix algebraic result: For every symmetric matrix  $\mathbf{B}$  there is an orthogonal transformation  $\mathbf{T}$  so that  $\mathbf{BT}^\top = \mathbf{T}^\top\mathbf{\Lambda}$  where  $\mathbf{\Lambda}$  is a diagonal matrix. Equivalently one could write  $\mathbf{B} = \mathbf{T}^\top\mathbf{\Lambda}\mathbf{T}$ . And if  $\mathbf{B}$  has rank  $r$ , then  $r$  of the diagonal elements are nonzero and the others zero. If one removes those eigenvectors from  $\mathbf{T}$  which belong to the eigenvalue zero, and calls the remaining matrix  $\mathbf{P}$ , one gets the following:

**THEOREM A.9.1.** *If  $\mathbf{B}$  is a symmetric  $n \times n$  matrix of rank  $r$ , then a  $r \times n$  matrix  $\mathbf{P}$  exists with  $\mathbf{PP}^\top = \mathbf{I}$  (any  $\mathbf{P}$  satisfying this condition which is not a square matrix is called incomplete orthogonal), and  $\mathbf{B} = \mathbf{P}^\top\mathbf{\Lambda}\mathbf{P}$ , where  $\mathbf{\Lambda}$  is a  $r \times r$  diagonal matrix with all diagonal elements nonzero.*

**PROOF.** Let  $\mathbf{T}$  be an orthogonal matrix whose rows are eigenvectors of  $\mathbf{B}$ , and partition it  $\mathbf{T} = \begin{bmatrix} \mathbf{P} \\ \mathbf{Q} \end{bmatrix}$  where  $\mathbf{P}$  consists of all eigenvectors with nonzero eigenvalue

(there are  $r$  of them). The eigenvalue property reads  $\mathbf{B} \begin{bmatrix} \mathbf{P}^\top & \mathbf{Q}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{P}^\top & \mathbf{Q}^\top \end{bmatrix} \begin{bmatrix} \mathbf{\Lambda} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}$

therefore by orthogonality  $\mathbf{T}^\top\mathbf{T} = \mathbf{I}$  follows  $\mathbf{B} = \begin{bmatrix} \mathbf{P}^\top & \mathbf{Q}^\top \end{bmatrix} \begin{bmatrix} \mathbf{\Lambda} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{P} \\ \mathbf{Q} \end{bmatrix} =$

$\mathbf{P}^\top\mathbf{\Lambda}\mathbf{P}$ . Orthogonality also means  $\mathbf{TT}^\top = \mathbf{I}$ , i.e.,  $\begin{bmatrix} \mathbf{P} \\ \mathbf{Q} \end{bmatrix} \begin{bmatrix} \mathbf{P}^\top & \mathbf{Q}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix}$ ,

therefore  $\mathbf{PP}^\top = \mathbf{I}$ .  $\square$

**PROBLEM 399.** *If  $\mathbf{B}$  is a  $n \times n$  symmetric matrix of rank  $r$  and  $\mathbf{B}^2 = \mathbf{B}$ , i.e.,  $\mathbf{B}$  is a projection, then a  $r \times n$  matrix  $\mathbf{P}$  exists with  $\mathbf{B} = \mathbf{P}^\top\mathbf{P}$  and  $\mathbf{PP}^\top = \mathbf{I}$ .*

**ANSWER.** Let  $\mathbf{t}$  be an eigenvector of the projection matrix  $\mathbf{B}$  with eigenvalue  $\lambda$ . Then  $\mathbf{B}^2\mathbf{t} = \mathbf{B}\mathbf{t}$ , i.e.,  $\lambda^2\mathbf{t} = \lambda\mathbf{t}$ , and since  $\mathbf{t} \neq \mathbf{o}$ ,  $\lambda^2 = \lambda$ . This is a quadratic equation with solutions  $\lambda = 0$  or  $\lambda = 1$ . The matrix  $\mathbf{\Lambda}$  from theorem A.9.1, whose diagonal elements are the nonzero eigenvalues, is therefore an identity matrix.  $\square$

A theorem similar to A.9.1 holds for arbitrary matrices. It is called the ‘singular value decomposition’:

**THEOREM A.9.2.** *Let  $\mathbf{B}$  be a  $m \times n$  matrix of rank  $r$ . Then  $\mathbf{B}$  can be expressed as*

$$(A.9.1) \quad \mathbf{B} = \mathbf{P}^\top\mathbf{\Lambda}\mathbf{Q}$$

where  $\mathbf{\Lambda}$  is a  $r \times r$  diagonal matrix with positive diagonal elements, and  $\mathbf{PP}^\top = \mathbf{I}$  as well as  $\mathbf{QQ}^\top = \mathbf{I}$ . The diagonal elements of  $\mathbf{\Lambda}$  are called the singular values of  $\mathbf{B}$ .

**PROOF.** If  $\mathbf{P}^\top\mathbf{\Lambda}\mathbf{Q}$  is the svd of  $\mathbf{B}$  then  $\mathbf{P}^\top\mathbf{\Lambda}\mathbf{Q}\mathbf{Q}^\top\mathbf{\Lambda}\mathbf{P} = \mathbf{P}^\top\mathbf{\Lambda}^2\mathbf{P}$  is the eigenvalue decomposition of  $\mathbf{BB}^\top$ . We will use this fact to construct  $\mathbf{P}$  and  $\mathbf{Q}$ , and then verify condition (A.9.1).  $\mathbf{P}$  and  $\mathbf{Q}$  have  $r$  rows each, write them

$$(A.9.2) \quad \mathbf{P} = \begin{bmatrix} \mathbf{p}_1^\top \\ \vdots \\ \mathbf{p}_r^\top \end{bmatrix} \quad \text{and} \quad \mathbf{Q} = \begin{bmatrix} \mathbf{q}_1^\top \\ \vdots \\ \mathbf{q}_r^\top \end{bmatrix}.$$

Then the  $\mathbf{p}_i$  are orthonormal eigenvectors of  $\mathbf{BB}^\top$  corresponding to the nonzero eigenvalues  $\lambda_i^2$ , and  $\mathbf{q}_i = \mathbf{B}^\top\mathbf{p}_i\lambda_i^{-1}$ . The proof that this definition is symmetric is left as exercise problem 400 below.

Now find  $\mathbf{p}_{r+1}, \dots, \mathbf{p}_m$  such that  $\mathbf{p}_1, \dots, \mathbf{p}_m$  is a complete set of orthonormal vectors, i.e.,  $\mathbf{p}_1\mathbf{p}_1^\top + \dots + \mathbf{p}_m\mathbf{p}_m^\top = \mathbf{I}$ . Then

$$(A.9.3) \quad \mathbf{B} = (\mathbf{p}_1\mathbf{p}_1^\top + \dots + \mathbf{p}_m\mathbf{p}_m^\top)\mathbf{B}$$

$$(A.9.4) \quad = (\mathbf{p}_1\mathbf{p}_1^\top + \dots + \mathbf{p}_r\mathbf{p}_r^\top)\mathbf{B} \quad \text{because } \mathbf{p}_i^\top\mathbf{B} = \mathbf{o}^\top \text{ for } i > r$$

$$(A.9.5) \quad = (\mathbf{p}_1\mathbf{q}_1^\top\lambda_1 + \dots + \mathbf{p}_r\mathbf{q}_r^\top\lambda_r) = \mathbf{P}^\top\mathbf{\Lambda}\mathbf{Q}.$$

**PROBLEM 400.** *Show that the  $\mathbf{q}_i$  are orthonormal eigenvectors of  $\mathbf{B}^\top\mathbf{B}$  corresponding to the same eigenvalues  $\lambda_i^2$ .*

**ANSWER.**

$$(A.9.6) \quad \mathbf{q}_i^\top\mathbf{q}_j = \lambda_i^{-1}\mathbf{p}_i^\top\mathbf{B}\mathbf{B}^\top\mathbf{p}_j\lambda_j^{-1} = \lambda_i^{-1}\mathbf{p}_i^\top\mathbf{p}_j\lambda_j^2\lambda_j^{-1} = \delta_{ij} \quad \text{Kronecker symbol}$$

$$(A.9.7) \quad \mathbf{B}^\top\mathbf{B}\mathbf{q}_i = \mathbf{B}^\top\mathbf{B}\mathbf{B}^\top\mathbf{p}_i\lambda_i^{-1} = \mathbf{B}^\top\mathbf{p}_i\lambda_i = \mathbf{q}_i\lambda_i^2$$

**PROBLEM 401.** *Show that  $\mathbf{B}\mathbf{q}_i = \lambda_i\mathbf{p}_i$  and  $\mathbf{B}^\top\mathbf{p}_i = \lambda_i\mathbf{q}_i$ .*

**ANSWER.** The second condition comes from the definition  $\mathbf{q}_i = \mathbf{B}^\top\mathbf{p}_i\lambda_i^{-1}$ , and premultiplied this definition by  $\mathbf{B}$  to get  $\mathbf{B}\mathbf{q}_i = \mathbf{B}\mathbf{B}^\top\mathbf{p}_i\lambda_i^{-1} = \lambda_i^2\mathbf{p}_i\lambda_i^{-1} = \lambda_i\mathbf{p}_i$ .

Let  $\mathbf{P}_0$  and  $\mathbf{Q}_0$  be such that  $\begin{bmatrix} \mathbf{P} \\ \mathbf{P}_0 \end{bmatrix}$  and  $\begin{bmatrix} \mathbf{Q} \\ \mathbf{Q}_0 \end{bmatrix}$  are orthogonal. Then the singular value decomposition can also be written in the full form, in which the matrix in the middle is  $m \times n$ :

$$(A.9.8) \quad \mathbf{B} = \begin{bmatrix} \mathbf{P}^\top & \mathbf{P}_0^\top \end{bmatrix} \begin{bmatrix} \mathbf{\Lambda} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{Q} \\ \mathbf{Q}_0 \end{bmatrix}$$

PROBLEM 402. Let  $\lambda_1$  be the biggest diagonal element of  $\mathbf{\Lambda}$ , and let  $\mathbf{c}$  and  $\mathbf{d}$  be two vectors with the properties that  $\mathbf{c}^\top \mathbf{B} \mathbf{d}$  is defined and  $\mathbf{c}^\top \mathbf{c} = 1$  as well as  $\mathbf{d}^\top \mathbf{d} = 1$ . Show that  $\mathbf{c}^\top \mathbf{B} \mathbf{d} \leq \lambda_1$ . The other singular values maximize among those who are orthogonal to the prior maximizers.

ANSWER.  $\mathbf{c}^\top \mathbf{B} \mathbf{d} = \mathbf{c}^\top \mathbf{P}^\top \mathbf{\Lambda} \mathbf{Q} \mathbf{d} = \mathbf{h}^\top \mathbf{\Lambda} \mathbf{k}$  where we call  $\mathbf{P} \mathbf{c} = \mathbf{h}$  and  $\mathbf{Q} \mathbf{d} = \mathbf{k}$ . By Cauchy-Schwartz (A.5.1),  $(\mathbf{h}^\top \mathbf{\Lambda} \mathbf{k})^2 \leq (\mathbf{h}^\top \mathbf{\Lambda} \mathbf{h})(\mathbf{k}^\top \mathbf{\Lambda} \mathbf{k})$ . Now  $(\mathbf{h}^\top \mathbf{\Lambda} \mathbf{k}) = \sum \lambda_{ii} h_i k_i \leq \sum \lambda_{11} h_i^2 = \lambda_{11} \mathbf{h}^\top \mathbf{h}$ . Now we only have to show that  $\mathbf{h}^\top \mathbf{h} \leq 1$ :  $1 - \mathbf{h}^\top \mathbf{h} = \mathbf{c}^\top \mathbf{c} - \mathbf{c}^\top \mathbf{P}^\top \mathbf{P} \mathbf{c} = \mathbf{c}^\top (\mathbf{I} - \mathbf{P}^\top \mathbf{P}) \mathbf{c} = \mathbf{c}^\top (\mathbf{I} - \mathbf{P}^\top \mathbf{P})(\mathbf{I} - \mathbf{P}^\top \mathbf{P}) \mathbf{c} \geq 0$ , here we used that  $\mathbf{P} \mathbf{P}^\top = \mathbf{I}$ , therefore  $\mathbf{P}^\top \mathbf{P}$  idempotent, therefore also  $\mathbf{I} - \mathbf{P}^\top \mathbf{P}$  idempotent.  $\square$

APPENDIX B

## Arrays of Higher Rank

This chapter was presented at the Array Programming Languages Conference in Berlin, on July 24, 2000.

Besides scalars, vectors, and matrices, also higher arrays are necessary in statistics; for instance, the “covariance matrix” of a random matrix is really an array of rank 4, etc. Usually, such higher arrays are avoided in the applied sciences because of the difficulties to write them on a two-dimensional sheet of paper. The following symbolic notation makes the structure of arrays explicit without writing them down element by element. It is hoped that this makes arrays easier to understand, and that this notation leads to simple high-level user interfaces for programming languages manipulating arrays.

### B.1. Informal Survey of the Notation

Each array is symbolized by a rectangular tile with arms sticking out, similar to a molecule. Tiles with one arm are vectors, those with two arms matrices, those with more arms are arrays of higher rank (or “valence” as in [SS35], [Mor73], and [MS86, p. 12]), and those without arms are scalars. The arrays considered here are rectangular, not “ragged,” therefore in addition to their rank we only need to know the dimension of each arm; it can be thought of as the number of fingers associated with this arm. Arrays can only hold hands (i.e., “contract” along two arms) if the hands have the same number of fingers.

Sometimes it is convenient to write the dimension of each arm at the end of the arm, i.e., a  $m \times n$  matrix  $\mathbf{A}$  can be represented as  $m$ — $\boxed{\mathbf{A}}$ — $n$ . Matrix products are represented by joining the obvious arms: if  $\mathbf{B}$  is  $n \times q$ , then the matrix product  $\mathbf{AB}$  is  $m$ — $\boxed{\mathbf{A}}$ — $n$ — $\boxed{\mathbf{B}}$ — $q$  or, in short,  $\boxed{\mathbf{A}}$ — $\boxed{\mathbf{B}}$ . The notation allows the reader to always tell which arm is which, even if the arms are not marked. If

$m$ — $\boxed{\mathbf{C}}$ — $r$  is  $m \times r$ , then the product  $\mathbf{C}^\top \mathbf{A}$  is

$$(B.1.1) \quad \mathbf{C}^\top \mathbf{A} = r \text{—} \boxed{\mathbf{C}} \text{—} m \text{—} \boxed{\mathbf{A}} \text{—} n = r \text{—} \boxed{\mathbf{C}} \text{—} m \text{—} \boxed{\mathbf{A}} \text{—} n .$$

In the second representation, the tile representing  $\mathbf{C}$  is turned by 180 degrees. Since the white part of the frame of  $\mathbf{C}$  is at the bottom, not on the top, one knows that the West arm of  $\mathbf{C}$ , not its East arm, is concatenated with the West arm of  $\mathbf{A}$ . The transpose of  $m$ — $\boxed{\mathbf{C}}$ — $r$  is  $r$ — $\boxed{\mathbf{C}}$ — $m$ , i.e., it is not a different entity but the same entity in a different position. The order in which the elements are arranged on the page (or in computer memory) is not a part of the definition of the array itself. Likewise, there is no distinction between row vectors and column vectors.

Vectors are usually, but not necessarily, written in such a way that their arms point West (column vector convention). If  $\boxed{\mathbf{a}}$  and  $\boxed{\mathbf{b}}$  are vectors, then their scalar product  $\mathbf{a}^\top \mathbf{b}$  is the concatenation  $\boxed{\mathbf{a}}$ — $\boxed{\mathbf{b}}$  which has no free arms, i.e., it is a scalar, and their outer product  $\mathbf{ab}^\top$  is  $\boxed{\mathbf{a}}$ — $\boxed{\mathbf{b}}$ , which is a matrix. Juxtaposition of tiles represents the outer product, i.e., the array consisting of the products of elements of the arrays represented by the tiles placed side by side.

The trace of a square matrix  $\boxed{\mathbf{Q}}$  is the concatenation  $\boxed{\mathbf{Q}}$ , which is a scalar since no arms are sticking out. In general, concatenation of two arms of the same tile represents *contraction*, i.e., summation over equal values of the indices associated with these two arms. This notation makes it obvious that  $\text{tr} \mathbf{XY} = \text{tr} \mathbf{YX}$ , because by definition there is no difference between  $\boxed{\mathbf{X}}$ — $\boxed{\mathbf{Y}}$  and

$$\boxed{\mathbf{Y}}$$
— $\boxed{\mathbf{X}}$ . Also  $\boxed{\mathbf{X}}$ — $\boxed{\mathbf{Y}}$  or  $\boxed{\mathbf{X}}$ — $\boxed{\mathbf{Y}}$  etc. represent the same array

(here array of rank zero, i.e., scalar). Each of these tiles can be evaluated in essentially two different ways. One way is

- (1) Juxtapose the tiles for  $\mathbf{X}$  and  $\mathbf{Y}$ , i.e., form their outer product, which is an array of rank 4 with typical element  $x_{mp}y_{qn}$ .
- (2) Connect the East arm of  $\mathbf{X}$  with the West arm of  $\mathbf{Y}$ . This is a contraction, resulting in an array of rank 2, the matrix product  $\mathbf{XY}$ , with typical element  $\sum_p x_{mp}y_{pn}$ .
- (3) Now connect the West arm of  $\mathbf{X}$  with the East arm of  $\mathbf{Y}$ . The result of this second contraction is a scalar, the trace  $\text{tr} \mathbf{XY} = \sum_{p,m} x_{mp}y_{pm}$ .

An alternative sequence of operations evaluating this same graph would be

- (1) Juxtapose the tiles for  $\mathbf{X}$  and  $\mathbf{Y}$ .
- (2) Connect the West arm of  $\mathbf{X}$  with the East arm of  $\mathbf{Y}$  to get the matrix product  $\mathbf{YX}$ .
- (3) Now connect the East arm of  $\mathbf{X}$  with the West arm of  $\mathbf{Y}$  to get  $\text{tr } \mathbf{YX}$ .

The result is the same, the notation does not specify which of these alternative evaluation paths is meant, and a computer receiving commands based on this notation can choose the most efficient evaluation path. Probably the most efficient evaluation path is given by (B.2.8) below: take the element-by-element product of  $\mathbf{X}$  with the transpose of  $\mathbf{Y}$ , and add all the elements of the resulting matrix.

If the user specifies  $\text{tr}(\mathbf{XY})$ , the computer is locked into one evaluation path: it first has to compute the matrix product  $\mathbf{XY}$ , even if  $\mathbf{X}$  is a column vector and  $\mathbf{Y}$  a row vector and it would be much more efficient to compute it as  $\text{tr}(\mathbf{YX})$ , and then form the trace, i.e., throw away all off-diagonal elements. If the trace is specified as  $\boxed{\mathbf{X}} \boxed{\mathbf{Y}}$ , the computer can choose the most efficient of a number of different evaluation paths transparently to the user. This advantage of the graphical notation is of course even more important if the graphs are more complex.

There is also the “diagonal” array, which in the case of rank 3 can be written

$$(B.1.2) \quad \begin{array}{c} n \\ \text{---} \end{array} \boxed{\Delta} \begin{array}{c} \text{---} \\ n \end{array} \quad \text{or} \quad \begin{array}{c} n \\ \text{---} \end{array} \boxed{\Delta} \begin{array}{c} \text{---} \\ n \end{array}$$

or similar configurations. It has 1’s down the main diagonal and 0’s elsewhere. It can be used to construct the diagonal matrix  $\text{diag}(\mathbf{x})$  of a vector (the square matrix with the vector in the diagonal and zeros elsewhere) as

$$(B.1.3) \quad \text{diag}(\mathbf{x}) = \begin{array}{c} n \\ \text{---} \end{array} \boxed{\Delta} \begin{array}{c} \text{---} \\ n \end{array} \boxed{\mathbf{x}}$$

the diagonal vector of a square matrix (i.e., the vector containing its diagonal elements) as

$$(B.1.4) \quad \text{---} \boxed{\Delta} \boxed{\mathbf{A}}$$

and the “Hadamard product” (element-by-element product) of two vectors  $\mathbf{x} * \mathbf{y}$

$$(B.1.5) \quad \mathbf{x} * \mathbf{y} = \text{---} \boxed{\Delta} \begin{array}{c} \boxed{\mathbf{x}} \\ \boxed{\mathbf{y}} \end{array}$$

All these are natural operations involving vectors and matrices, but the usual matrix notation cannot represent them and therefore ad-hoc notation must be invented for them. In our graphical representation, however, they all can be built up from a small number of atomic operations, which will be enumerated in Section B.2.

Each such graph can be evaluated in a number of different ways, and all the evaluations give the same result. In principle, each graph can be evaluated as follows: form the outer product of all arrays involved, and then contract along all those pairs of arms which are connected. For practical implementations it is more efficient to develop functions which connect two arrays along one or several of their arms without first forming outer products, and to perform the array concatenations recursively in such a way that contractions are done as early as possible. A computer might be programmed to decide on the most efficient construction path for any given array.

### B.2. Axiomatic Development of Array Operations

The following sketch shows how this axiom system might be built up. Since I am an economist I do not plan to develop the material presented here any further. Others are invited to take over. If you are interested in working on this, I would be happy to hear from you; email me at [ehrbbar@econ.utah.edu](mailto:ehrbbar@econ.utah.edu)

There are two kinds of special arrays: unit vectors and diagonal arrays.

For every natural number  $m \geq 1$ ,  $m$  unit vectors  $m \text{---} \boxed{\mathbf{i}}$  ( $i = 1, \dots, m$ ) exist. Despite the fact that the unit vectors are denoted here by numbers, there is an intrinsic ordering among them; they might as well have the names “red, green, blue, ...” (From (B.2.4) and other axioms below it will follow that each unit vector can be represented as a  $m$ -vector with 1 as one of the components and 0 elsewhere.)

For every rank  $\geq 1$  and dimension  $n \geq 1$  there is a unique diagonal array denoted by  $\Delta$ . Their main properties are (B.2.1) and (B.2.2). (This and the other axioms must be formulated in such a way that it will be possible to show that the diagonal arrays of rank 1 are the “vectors of ones”  $\mathbf{1}$  which have 1 in every component; diagonal arrays of rank 2 are the identity matrices; and for higher ranks, all arms of a diagonal array have the same dimension, and their  $ijk \dots$  element is 1 if  $i = j = k = \dots$  and 0 otherwise.) Perhaps it makes sense to define the diagonal array of rank 1 and dimension  $n$  to be the scalar  $n$ , and to declare all arrays which are everywhere 0-dimensional to be diagonal.



There are only three operations of arrays: their outer product, represented by writing them side by side, contraction, represented by the joining of arms, and the direct sum, which will be defined now:

The direct sum is the operation by which a vector can be built up from scalars, a matrix from its row or column vectors, an array of rank 3 from its layers, etc. The direct sum of a set of  $r$  similar arrays (i.e., arrays which have the same number of arms, and corresponding arms have the same dimensions) is an array which has one additional arm, called the reference arm of the direct sum. If one “saturates” the reference arm with the  $i$ th unit vector, one gets the  $i$ th original array back, and this property defines the direct sum uniquely:

$$\bigoplus_{i=1}^r \begin{array}{c} m \\ \boxed{A_i} \\ q \end{array} \text{---} n = r \text{---} \begin{array}{c} m \\ \boxed{S} \\ q \end{array} \text{---} n \Rightarrow \boxed{i} \text{---} r \text{---} \begin{array}{c} m \\ \boxed{S} \\ q \end{array} \text{---} n = \begin{array}{c} m \\ \boxed{A_i} \\ q \end{array} \text{---} n .$$

It is impossible to tell which is the first summand and which the second, direct sum is an operation defined on finite sets of arrays (where different elements of a set may be equal to each other in every respect but still have different identities).

There is a broad rule of associativity: the order in which outer products and contractions are performed does not matter, as long as the at the end, the right arms are connected with each other. And there are distributive rules involving (contracted) outer products and direct sums.

Additional rules apply for the special arrays. If two different diagonal arrays join arms, the result is again a diagonal array. For instance, the following three concatenations of diagonal three-way arrays are identical, and they all evaluate to the (for a given dimension) unique diagonal array or rank 4:

$$(B.2.1) \quad \begin{array}{c} \diagup \quad \diagdown \\ \boxed{\Delta} \\ \diagdown \quad \diagup \\ \boxed{\Delta} \\ \diagup \quad \diagdown \end{array} = \begin{array}{c} \diagup \quad \diagdown \\ \boxed{\Delta} \text{---} \boxed{\Delta} \\ \diagdown \quad \diagup \end{array} = \begin{array}{c} \diagup \quad \diagdown \\ \boxed{\Delta} \\ \diagdown \quad \diagup \\ \boxed{\Delta} \\ \diagup \quad \diagdown \end{array} = \begin{array}{c} \diagup \quad \diagdown \\ \boxed{\Delta} \\ \diagdown \quad \diagup \end{array}$$

The diagonal array of rank 2 is neutral under concatenation, i.e., it can be written as

$$(B.2.2) \quad n \text{---} \boxed{\Delta} \text{---} n = \text{---} .$$

because attaching it to any array will not change this array. (B.2.1) and (B.2.2) make it possible to represent diagonal arrays simply as the branching points of several arms.

This will make the array notation even simpler. However in the present introductory article, all diagonal arrays will be shown explicitly, and the vector of ones will be denoted  $m \text{---} \boxed{\iota} \text{---} m$  instead of  $m \text{---} \boxed{\Delta} \text{---} m$  or perhaps  $m \text{---} \boxed{\delta} \text{---} m$ .

Unit vectors concatenate as follows:

$$(B.2.3) \quad \boxed{i} \text{---} m \text{---} \boxed{j} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

and the direct sum of all unit vectors is the diagonal array of rank 2:

$$(B.2.4) \quad \bigoplus_{i=1}^n \boxed{i} \text{---} n = n \text{---} \boxed{\Delta} \text{---} n = \text{---} .$$

I am sure there will be modifications if one works it all out in detail, but if done right, the number of axioms should be fairly small. Element-by-element addition of arrays is not an axiom because it can be derived: if one saturates the reference arm of a direct sum with the vector of ones, one gets the element-by-element sum of the arrays in this direct sum. Multiplication of an array by a scalar is also contained in the above system of axioms: it is simply the outer product with an array of rank zero.

**PROBLEM 403.** Show that the saturation of an arm of a diagonal array with the vector of ones is the same as dropping this arm.

**ANSWER.** Since the vector of ones is the diagonal array of rank 1, this is a special case of the general concatenation rule for diagonal arrays.

**PROBLEM 404.** Show that the diagonal matrix of the vector of ones is the identity matrix, i.e.,

$$(B.2.5) \quad n \text{---} \boxed{\Delta} \text{---} n \text{---} \boxed{\iota} \text{---} n = \text{---} .$$

**ANSWER.** In view of (B.2.2), this is a special case of Problem 403.

**PROBLEM 405.** A trivial array operation is the addition of an arm of dimension 1; for instance, this is how a  $n$ -vector can be turned into a  $n \times 1$  matrix. Is this operation contained in the above system of axioms?

**ANSWER.** It is a special case of the direct sum: the direct sum of one array only, the only effect of which is the addition of the reference arm.

From (B.2.4) and (B.2.2) follows that every array of rank  $k$  can be represented as a direct sum of arrays of rank  $k - 1$ , and recursively, as iterated direct sums of those scalars which one gets by saturating all arms with unit vectors. Hence t

following “extensionality property”: if the arrays  $A$  and  $B$  are such that for all possible conformable choices of unit vectors  $\kappa_1 \cdots \kappa_8$  follows

$$(B.2.6) \quad \begin{array}{c} \boxed{\kappa_3} \\ \boxed{\kappa_2} \\ \boxed{\kappa_1} \end{array} \begin{array}{c} \boxed{\kappa_4} \\ \boxed{A} \\ \boxed{\kappa_8} \end{array} \begin{array}{c} \boxed{\kappa_5} \\ \boxed{\kappa_6} \\ \boxed{\kappa_7} \end{array} = \begin{array}{c} \boxed{\kappa_3} \\ \boxed{\kappa_2} \\ \boxed{\kappa_1} \end{array} \begin{array}{c} \boxed{\kappa_4} \\ \boxed{B} \\ \boxed{\kappa_8} \end{array} \begin{array}{c} \boxed{\kappa_5} \\ \boxed{\kappa_6} \\ \boxed{\kappa_7} \end{array}$$

then  $A = B$ . This is why the saturation of an array with unit vectors can be considered one of its “elements,” i.e.,

$$(B.2.7) \quad \begin{array}{c} \boxed{\kappa_3} \\ \boxed{\kappa_2} \\ \boxed{\kappa_1} \end{array} \begin{array}{c} \boxed{\kappa_4} \\ \boxed{A} \\ \boxed{\kappa_8} \end{array} \begin{array}{c} \boxed{\kappa_5} \\ \boxed{\kappa_6} \\ \boxed{\kappa_7} \end{array} = a_{\kappa_1 \kappa_2 \kappa_3 \kappa_4 \kappa_5 \kappa_6 \kappa_7 \kappa_8}$$

From (B.2.3) and (B.2.4) follows that the concatenation of two arrays by joining one or more pairs of arms consists in forming all possible products and summing over those subscripts (arms) which are joined to each other. For instance, if

$$m \text{---} \boxed{A} \text{---} n \text{---} \boxed{B} \text{---} r = m \text{---} \boxed{C} \text{---} r,$$

then  $c_{\mu\rho} = \sum_{\nu=1}^n a_{\mu\nu} b_{\nu\rho}$ . This is one of the most basic facts if one thinks of arrays as collections of elements. From this point of view, the proposed notation is simply a graphical elaboration of Einstein’s summation convention. But in the holistic approach taken by the proposed system of axioms, which is informed by category theory, it is an implication; it comes at the end, not the beginning.

Instead of considering arrays as bags filled with elements, with the associated false problem of specifying the order in which the elements are packed into the bag, this notation and system of axioms consider each array as an abstract entity, associated with a certain finite graph. These entities can be operated on as specified in the axioms, but the only time they lose their abstract character is when they are fully saturated, i.e., concatenated with each other in such a way that no free arms are left: in this case they become scalars. An array of rank 1 is not the same as a vector, although it can be *represented* as a vector—after an ordering of its elements has been specified. This ordering is not part of the definition of the array itself. (Some vectors, such as time series, have an intrinsic ordering, but I am speaking here of the simplest case where they do not.) Also the ordering of the arms is not specified, and the order in which a set of arrays is packed into its direct sum is not specified

either. These axioms therefore make a strict distinction between the abstract entities themselves (which the user is interested in) and their various representations (which the computer worries about).

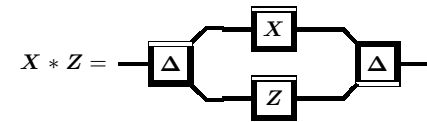
Maybe the following examples may clarify these points. If you specify a set of colors as {red, green, blue}, then this representation has an ordering built in: red comes first, then green, then blue. However this ordering is not part of the definition of the set; {green, red, blue} is the same set. The two notations are two different representations of the same set. Another example: mathematicians usually distinguish between the outer products  $A \otimes B$  and  $B \otimes A$ ; there is a “natural isomorphism” between them but they are two different objects. In the system of axioms proposed here these two notations are two different representations of the same object, as in the set example. This object is represented by a graph which has  $A$  and  $B$  as nodes, but it is not apparent from this graph which node comes first. Interesting conceptual issues are involved here. The proposed axioms are quite different than e.g. [Mor7].

PROBLEM 406. The trace of the product of two matrices can be written as

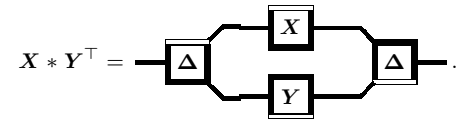
$$(B.2.8) \quad \text{tr}(XY) = \iota^\top (X * Y^\top) \iota.$$

I.e., one forms the element-by-element product of  $X$  and  $Y^\top$  and takes the sum of all the elements of the resulting matrix. Use tile notation to show that this gives indeed  $\text{tr}(XY)$ .

ANSWER. In analogy with (B.1.5), the Hadamard product of the two matrices  $X$  and  $Z$ , is their element by element multiplication, is



If  $Z = Y^\top$ , one gets

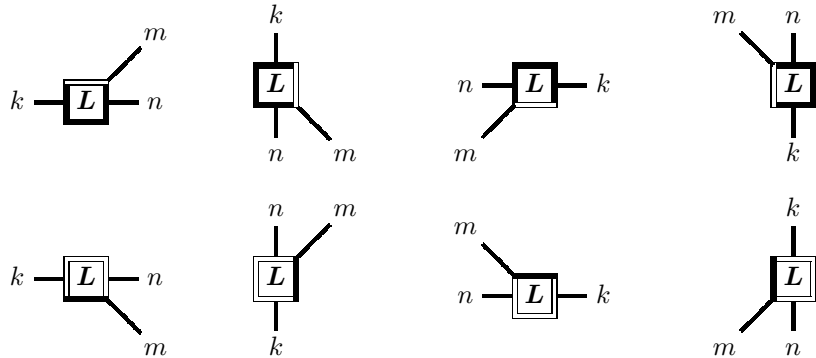


Therefore one gets, using (B.2.5):

$$\iota^\top (X * Y^\top) \iota = \begin{array}{c} \boxed{\iota} \\ \boxed{\Delta} \end{array} \begin{array}{c} \boxed{X} \\ \boxed{Y} \end{array} \begin{array}{c} \boxed{\Delta} \\ \boxed{\iota} \end{array} = \begin{array}{c} \boxed{X} \\ \boxed{Y} \end{array} = \text{tr}(XY)$$

### B.3. An Additional Notational Detail

Besides turning a tile by 90, 180, or 270 degrees, the notation proposed here also allows to flip the tile over. The tile  $\square$  (here drawn without its arms) is simply the tile  $\square$  laid on its face; i.e., those parts of the frame, which are black on the side visible to the reader, are white on the opposite side and vice versa. If one flips a tile, the arms appear in a mirror-symmetric manner. For a matrix, flipping over is equivalent to turning by 180 degrees, i.e., there is no difference between the matrix  $\square$  and the matrix  $\square$ . Since sometimes one and sometimes the other notation seems more natural, both will be used. For higher arrays, flipping over arranges the arms in a different fashion, which is sometimes convenient in order to keep the graphs uncluttered. It will be especially useful for differentiation. If one allows turning in 90 degree increments and flipping, each array can be represented in eight different positions, as shown here with a hypothetical array of rank 3:



The black-and-white pattern at the edge of the tile indicates whether and how much the tile has been turned and/or flipped over, so that one can keep track which arm is which. In the above example, the arm with dimension  $k$  will always be called the West arm, whatever position the tile is in.

### B.4. Equality of Arrays and Extended Substitution

Given the flexibility of representing the same array in various positions for concatenation, specific conventions are necessary to determine when two such arrays in generalized positions are equal to each other. Expressions like

$$\square = \square \quad \text{or} \quad \square = \square$$

are not allowed. The arms on both sides of the equal sign must be parallel, in order to make it clear which arm corresponds to which. A permissible way to write the above expressions would therefore be

$$\square = \square \quad \text{and} \quad \square = \square$$

One additional benefit of this tile notation is the ability to substitute arrays with different numbers of arms into an equation. This is also a necessity since the number of possible arms is unbounded. This multiplicity can only be coped with because each arm in an identity written in this notation can be replaced by a bundle of many arms.

Extended substitution also makes it possible to extend definitions familiar from matrices to higher arrays. For instance we want to be able to say that the array  $\square$  is symmetric if and only if  $\square = \square$ . This notion of symmetry is not limited to arrays of rank 2. The arms of this array may symbolize not just

a single arm, but whole bundles of arms; for instance an array of the form  $\square$

satisfying  $\square = \square$  is symmetric according to this definition, and so

every scalar. Also the notion of a nonnegative definite matrix, or of a matrix inverse or generalized inverse, or of a projection matrix, can be extended to arrays in this way.

### B.5. Vectorization and Kronecker Product

One conventional generally accepted method to deal with arrays of rank  $> 2$  is the Kronecker product. If  $A$  and  $B$  are both matrices, then the outer product in this notation is

$$(B.5.1) \quad \square$$

Since this is an array of rank 4, there is no natural way to write its elements down on a sheet of paper. This is where the Kronecker product steps in. The Kronecker product of two matrices is their outer product written again as a matrix. Its definition includes a protocol how to arrange the elements of an array of rank 4 as a matrix. Alongside the Kronecker product, also the vectorization operator is useful, which

a protocol how to arrange the elements of a matrix as a vector, and also the so-called “commutation matrices” may become necessary. Here are the relevant definitions:

**B.5.1. Vectorization of a Matrix.** If  $\mathbf{A}$  is a matrix, then  $\text{vec}(\mathbf{A})$  is the vector obtained by stacking the column vectors on top of each other, i.e.,

$$(B.5.2) \quad \text{if } \mathbf{A} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_n] \quad \text{then} \quad \text{vec}(\mathbf{A}) = \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{bmatrix}.$$

The vectorization of a matrix is merely a different arrangement of the elements of the matrix on paper, just as the transpose of a matrix.

PROBLEM 407. Show that  $\text{tr}(\mathbf{B}^\top \mathbf{C}) = (\text{vec } \mathbf{B})^\top \text{vec } \mathbf{C}$ .

ANSWER. Both sides are  $\sum b_{ji}c_{ji}$ . (B.5.28) is a proof in tile notation which does not have to look at the matrices involved element by element.  $\square$

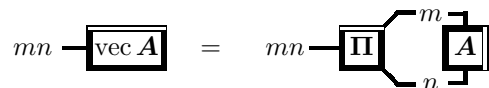
By the way, a better protocol for vectorizing would have been to assemble all rows into one long row vector and then converting it into a column vector. In other words

$$\text{if } \mathbf{B} = \begin{bmatrix} \mathbf{b}_1^\top \\ \vdots \\ \mathbf{b}_m^\top \end{bmatrix} \quad \text{then } \text{vec}(\mathbf{B}) \text{ should have been defined as } \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_m \end{bmatrix}.$$

The usual protocol of stacking the columns is inconsistent with the lexicographical ordering used in the Kronecker product. Using the alternative definition, equation (B.5.19) which will be discussed below would be a little more intelligible; it would read

$$\text{vec}(\mathbf{ABC}) = (\mathbf{A} \otimes \mathbf{C}^\top) \text{vec } \mathbf{B} \quad \text{with the alternative definition of } \text{vec}$$

and also the definition of vectorization in tile notation would be a little less awkward; instead of (B.5.24) one would have



But this is merely a side remark; we will use the conventional definition (B.5.2) throughout.

**B.5.2. Kronecker Product of Matrices.** Let  $\mathbf{A}$  and  $\mathbf{B}$  be two matrices,  $\mathbf{A}$  is  $m \times n$  and  $\mathbf{B}$  is  $r \times q$ . Their Kronecker product  $\mathbf{A} \otimes \mathbf{B}$  is the  $mr \times nq$  matrix which in partitioned form can be written

$$(B.5.3) \quad \mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}$$

This convention of how to write the elements of an array of rank 4 as a matrix is not symmetric, so that usually  $\mathbf{A} \otimes \mathbf{C} \neq \mathbf{C} \otimes \mathbf{A}$ . Both Kronecker products represent the same abstract array, but they arrange it differently on the page. However, in many other respects, the Kronecker product maintains the properties of outer products

PROBLEM 408. [The71, pp. 303–306] Prove the following simple properties of the Kronecker product:

$$(B.5.4) \quad (\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top$$

$$(B.5.5) \quad (\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C})$$

$$(B.5.6) \quad \mathbf{I} \otimes \mathbf{I} = \mathbf{I}$$

$$(B.5.7) \quad (\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}$$

$$(B.5.8) \quad (\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$$

$$(B.5.9) \quad (\mathbf{A} \otimes \mathbf{B})^- = \mathbf{A}^- \otimes \mathbf{B}^-$$

$$(B.5.10) \quad \mathbf{A} \otimes (\mathbf{B} + \mathbf{C}) = \mathbf{A} \otimes \mathbf{B} + \mathbf{A} \otimes \mathbf{C}$$

$$(B.5.11) \quad (\mathbf{A} + \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes \mathbf{C} + \mathbf{B} \otimes \mathbf{C}$$

$$(B.5.12) \quad (c\mathbf{A}) \otimes \mathbf{B} = \mathbf{A} \otimes (c\mathbf{B}) = c(\mathbf{A} \otimes \mathbf{B})$$

$$(B.5.13) \quad \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \otimes \mathbf{B} = \begin{bmatrix} \mathbf{A}_{11} \otimes \mathbf{B} & \mathbf{A}_{12} \otimes \mathbf{B} \\ \mathbf{A}_{21} \otimes \mathbf{B} & \mathbf{A}_{22} \otimes \mathbf{B} \end{bmatrix}$$

$$(B.5.14) \quad \text{rank}(\mathbf{A} \otimes \mathbf{B}) = (\text{rank } \mathbf{A})(\text{rank } \mathbf{B})$$

$$(B.5.15) \quad \text{tr}(\mathbf{A} \otimes \mathbf{B}) = (\text{tr } \mathbf{A})(\text{tr } \mathbf{B})$$

If  $\mathbf{a}$  is a  $1 \times 1$  matrix, then

$$(B.5.16) \quad \mathbf{a} \otimes \mathbf{B} = \mathbf{B} \otimes \mathbf{a} = \mathbf{aB}$$

$$(B.5.17) \quad \det(\mathbf{A} \otimes \mathbf{B}) = (\det(\mathbf{A}))^n (\det(\mathbf{B}))^k$$

where  $\mathbf{A}$  is  $k \times k$  and  $\mathbf{B}$  is  $n \times n$ .

ANSWER. For the determinant use the following facts: if  $\mathbf{a}$  is an eigenvector of  $\mathbf{A}$  with eigenvalue  $\alpha$  and  $\mathbf{b}$  is an eigenvector of  $\mathbf{B}$  with eigenvalue  $\beta$ , then  $\mathbf{a} \otimes \mathbf{b}$  is an eigenvector of  $\mathbf{A} \otimes \mathbf{B}$  with eigenvalue  $\alpha\beta$ . The determinant is the product of all eigenvalues (multiple eigenvalues being counted several times). Count how many there are.

An alternative approach would be to write  $\mathbf{A} \otimes \mathbf{B} = (\mathbf{A} \otimes \mathbf{I})(\mathbf{I} \otimes \mathbf{B})$  and then to argue that  $\det(\mathbf{A} \otimes \mathbf{I}) = (\det(\mathbf{A}))^n$  and  $\det(\mathbf{I} \otimes \mathbf{B}) = (\det(\mathbf{B}))^k$ .

The formula for the rank can be shown using  $\text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A}\mathbf{A}^\top)$ . compare Problem 366.  $\square$

PROBLEM 409. 2 points [JHG<sup>+</sup>88, pp. 962–4] Write down the Kronecker product of

$$(B.5.18) \quad \mathbf{A} = \begin{bmatrix} 1 & 3 \\ 2 & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} 2 & 2 & 0 \\ 1 & 0 & 3 \end{bmatrix}.$$

Show that  $\mathbf{A} \otimes \mathbf{B} \neq \mathbf{B} \otimes \mathbf{A}$ . Which other facts about the outer product do not carry over to the Kronecker product?

ANSWER.

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} 2 & 2 & 0 & 6 & 6 & 0 \\ 1 & 0 & 3 & 3 & 0 & 9 \\ 4 & 4 & 0 & 0 & 0 & 0 \\ 2 & 0 & 6 & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{B} \otimes \mathbf{A} = \begin{bmatrix} 2 & 6 & 2 & 6 & 0 & 0 \\ 4 & 0 & 4 & 0 & 0 & 0 \\ 1 & 3 & 0 & 0 & 3 & 9 \\ 2 & 0 & 0 & 0 & 6 & 0 \end{bmatrix}$$

Partitioning of the matrix on the right does not carry over.  $\square$

PROBLEM 410. [JHG<sup>+</sup>88, p. 965] Show that

$$(B.5.19) \quad \text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{B}).$$

ANSWER. Assume  $\mathbf{A}$  is  $k \times m$ ,  $\mathbf{B}$  is  $m \times n$ , and  $\mathbf{C}$  is  $n \times p$ . Write  $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_k^\top \end{bmatrix}$  and  $\mathbf{B} =$

$\begin{bmatrix} \mathbf{b}_1 & \cdots & \mathbf{b}_n \end{bmatrix}$ . Then  $(\mathbf{C}^\top \otimes \mathbf{A}) \text{vec} \mathbf{B} =$

$$= \begin{bmatrix} c_{11}\mathbf{A} & c_{21}\mathbf{A} & \cdots & c_{n1}\mathbf{A} \\ c_{12}\mathbf{A} & c_{22}\mathbf{A} & \cdots & c_{n2}\mathbf{A} \\ \vdots & \vdots & \ddots & \vdots \\ c_{1p}\mathbf{A} & c_{2p}\mathbf{A} & \cdots & c_{np}\mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_n \end{bmatrix} = \begin{bmatrix} c_{11}\mathbf{a}_1^\top \mathbf{b}_1 + c_{21}\mathbf{a}_1^\top \mathbf{b}_2 + \cdots + c_{n1}\mathbf{a}_1^\top \mathbf{b}_n \\ c_{11}\mathbf{a}_2^\top \mathbf{b}_1 + c_{21}\mathbf{a}_2^\top \mathbf{b}_2 + \cdots + c_{n1}\mathbf{a}_2^\top \mathbf{b}_n \\ \vdots \\ c_{11}\mathbf{a}_k^\top \mathbf{b}_1 + c_{21}\mathbf{a}_k^\top \mathbf{b}_2 + \cdots + c_{n1}\mathbf{a}_k^\top \mathbf{b}_n \\ c_{12}\mathbf{a}_1^\top \mathbf{b}_1 + c_{22}\mathbf{a}_1^\top \mathbf{b}_2 + \cdots + c_{n2}\mathbf{a}_1^\top \mathbf{b}_n \\ c_{12}\mathbf{a}_2^\top \mathbf{b}_1 + c_{22}\mathbf{a}_2^\top \mathbf{b}_2 + \cdots + c_{n2}\mathbf{a}_2^\top \mathbf{b}_n \\ \vdots \\ c_{12}\mathbf{a}_k^\top \mathbf{b}_1 + c_{22}\mathbf{a}_k^\top \mathbf{b}_2 + \cdots + c_{n2}\mathbf{a}_k^\top \mathbf{b}_n \\ \vdots \\ c_{1p}\mathbf{a}_1^\top \mathbf{b}_1 + c_{2p}\mathbf{a}_1^\top \mathbf{b}_2 + \cdots + c_{np}\mathbf{a}_1^\top \mathbf{b}_n \\ c_{1p}\mathbf{a}_2^\top \mathbf{b}_1 + c_{2p}\mathbf{a}_2^\top \mathbf{b}_2 + \cdots + c_{np}\mathbf{a}_2^\top \mathbf{b}_n \\ \vdots \\ c_{1p}\mathbf{a}_k^\top \mathbf{b}_1 + c_{2p}\mathbf{a}_k^\top \mathbf{b}_2 + \cdots + c_{np}\mathbf{a}_k^\top \mathbf{b}_n \end{bmatrix}.$$

One obtains the same result by vectorizing the matrix

$$\begin{aligned} \mathbf{ABC} &= \begin{bmatrix} \mathbf{a}_1^\top \mathbf{b}_1 & \mathbf{a}_1^\top \mathbf{b}_2 & \cdots & \mathbf{a}_1^\top \mathbf{b}_n \\ \mathbf{a}_2^\top \mathbf{b}_1 & \mathbf{a}_2^\top \mathbf{b}_2 & \cdots & \mathbf{a}_2^\top \mathbf{b}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_k^\top \mathbf{b}_1 & \mathbf{a}_k^\top \mathbf{b}_2 & \cdots & \mathbf{a}_k^\top \mathbf{b}_n \end{bmatrix} \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1p} \\ c_{21} & c_{22} & \cdots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \cdots & c_{np} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{a}_1^\top \mathbf{b}_1 c_{11} + \mathbf{a}_1^\top \mathbf{b}_2 c_{21} + \cdots + \mathbf{a}_1^\top \mathbf{b}_n c_{n1} & \mathbf{a}_1^\top \mathbf{b}_1 c_{12} + \mathbf{a}_1^\top \mathbf{b}_2 c_{22} + \cdots + \mathbf{a}_1^\top \mathbf{b}_n c_{n2} & \cdots \\ \mathbf{a}_2^\top \mathbf{b}_1 c_{11} + \mathbf{a}_2^\top \mathbf{b}_2 c_{21} + \cdots + \mathbf{a}_2^\top \mathbf{b}_n c_{n1} & \mathbf{a}_2^\top \mathbf{b}_1 c_{12} + \mathbf{a}_2^\top \mathbf{b}_2 c_{22} + \cdots + \mathbf{a}_2^\top \mathbf{b}_n c_{n2} & \cdots \\ \vdots & \vdots & \ddots \\ \mathbf{a}_k^\top \mathbf{b}_1 c_{11} + \mathbf{a}_k^\top \mathbf{b}_2 c_{21} + \cdots + \mathbf{a}_k^\top \mathbf{b}_n c_{n1} & \mathbf{a}_k^\top \mathbf{b}_1 c_{12} + \mathbf{a}_k^\top \mathbf{b}_2 c_{22} + \cdots + \mathbf{a}_k^\top \mathbf{b}_n c_{n2} & \cdots \\ \cdots & \mathbf{a}_1^\top \mathbf{b}_1 c_{1p} + \mathbf{a}_1^\top \mathbf{b}_2 c_{2p} + \cdots + \mathbf{a}_1^\top \mathbf{b}_n c_{np} \\ \cdots & \mathbf{a}_2^\top \mathbf{b}_1 c_{1p} + \mathbf{a}_2^\top \mathbf{b}_2 c_{2p} + \cdots + \mathbf{a}_2^\top \mathbf{b}_n c_{np} \\ \vdots & \vdots \\ \cdots & \mathbf{a}_k^\top \mathbf{b}_1 c_{1p} + \mathbf{a}_k^\top \mathbf{b}_2 c_{2p} + \cdots + \mathbf{a}_k^\top \mathbf{b}_n c_{np} \end{bmatrix} \end{aligned}$$

The main challenge in this automatic proof is to fit the many matrix rows, columns, and scalar elements involved on the same sheet of paper. Among the shuffling of matrix entries, it is easy to lose track of how the result comes about. Later, in equation (B.5.29), a compact and intelligible proof will be given in tile notation.

The dispersion of a random matrix  $\mathbf{Y}$  is often given as the matrix  $\mathcal{V}[\text{vec} \mathbf{Y}]$ , where the vectorization is usually not made explicit, i.e., this matrix is denoted  $\mathcal{V}[\mathbf{Y}]$ .

PROBLEM 411. If  $\mathcal{V}[\text{vec} \mathbf{Y}] = \mathbf{\Sigma} \otimes \mathbf{\Omega}$  and  $\mathbf{P}$  and  $\mathbf{Q}$  are matrices of constant rank, show that  $\mathcal{V}[\text{vec} \mathbf{PYQ}] = (\mathbf{Q}^\top \mathbf{\Sigma} \mathbf{Q}) \otimes (\mathbf{P} \mathbf{\Omega} \mathbf{P}^\top)$ .

ANSWER. Apply (B.5.19):  $\mathcal{V}[\text{vec} \mathbf{PYQ}] = \mathcal{V}[(\mathbf{Q}^\top \otimes \mathbf{P}) \text{vec} \mathbf{Y}] = (\mathbf{Q}^\top \otimes \mathbf{P})(\mathbf{\Sigma} \otimes \mathbf{\Omega})(\mathbf{Q} \otimes \mathbf{P})$ . Now apply (B.5.7).

PROBLEM 412. 2 points If  $\alpha$  and  $\gamma$  are vectors, then show that  $\text{vec}(\alpha \gamma^\top) = \gamma \otimes \alpha$ .

ANSWER. One sees this by writing down the matrices, or one can use (B.5.19) with  $\mathbf{A} = \mathbf{B} = \mathbf{1}$ , the  $1 \times 1$  matrix, and  $\mathbf{C} = \gamma^\top$ .

PROBLEM 413. 2 points If  $\alpha$  is a nonrandom vector and  $\delta$  a random vector, show that  $\mathcal{V}[\delta \otimes \alpha] = \mathcal{V}[\delta] \otimes (\alpha \alpha^\top)$ .

ANSWER.

$$\delta \otimes \alpha = \begin{bmatrix} \alpha \delta_1 \\ \vdots \\ \alpha \delta_n \end{bmatrix} \quad \mathcal{V}[\delta \otimes \alpha] = \begin{bmatrix} \alpha \operatorname{var}[\delta_1] \alpha^\top & \alpha \operatorname{cov}[\delta_1, \delta_2] \alpha^\top & \cdots & \alpha \operatorname{cov}[\delta_1, \delta_n] \alpha^\top \\ \alpha \operatorname{cov}[\delta_2, \delta_1] \alpha^\top & \alpha \operatorname{var}[\delta_2] \alpha^\top & \cdots & \alpha \operatorname{cov}[\delta_2, \delta_n] \alpha^\top \\ \vdots & \vdots & \ddots & \vdots \\ \alpha \operatorname{cov}[\delta_n, \delta_1] \alpha^\top & \alpha \operatorname{cov}[\delta_n, \delta_2] \alpha^\top & \cdots & \alpha \operatorname{cov}[\delta_n, \delta_n] \alpha^\top \end{bmatrix} =$$

$$= \begin{bmatrix} \operatorname{var}[\delta_1] \alpha \alpha^\top & \operatorname{cov}[\delta_1, \delta_2] \alpha \alpha^\top & \cdots & \operatorname{cov}[\delta_1, \delta_n] \alpha \alpha^\top \\ \operatorname{cov}[\delta_2, \delta_1] \alpha \alpha^\top & \operatorname{var}[\delta_2] \alpha \alpha^\top & \cdots & \operatorname{cov}[\delta_2, \delta_n] \alpha \alpha^\top \\ \vdots & \vdots & \ddots & \vdots \\ \operatorname{cov}[\delta_n, \delta_1] \alpha \alpha^\top & \operatorname{cov}[\delta_n, \delta_2] \alpha \alpha^\top & \cdots & \operatorname{cov}[\delta_n, \delta_n] \alpha \alpha^\top \end{bmatrix} = \mathcal{V}[\delta] \otimes \alpha \alpha^\top$$

□

**B.5.3. The Commutation Matrix.** Besides the Kronecker product and the vectorization operator, also the “commutation matrix” [MN88, pp. 46/7], [Mag88, p. 35] is needed for certain operations involving arrays of higher rank. Assume  $\mathbf{A}$  is  $m \times n$ . Then the commutation matrix  $\mathbf{K}^{(m,n)}$  is the  $mn \times mn$  matrix which transforms  $\operatorname{vec} \mathbf{A}$  into  $\operatorname{vec}(\mathbf{A}^\top)$ :

$$(B.5.20) \quad \mathbf{K}^{(m,n)} \operatorname{vec} \mathbf{A} = \operatorname{vec}(\mathbf{A}^\top)$$

The main property of the commutation matrix is that it allows to commute the Kronecker product. For any  $m \times n$  matrix  $\mathbf{A}$  and  $r \times q$  matrix  $\mathbf{B}$  follows

$$(B.5.21) \quad \mathbf{K}^{(r,m)} (\mathbf{A} \otimes \mathbf{B}) \mathbf{K}^{(n,q)} = \mathbf{B} \otimes \mathbf{A}$$

PROBLEM 414. Use (B.5.20) to compute  $\mathbf{K}^{(2,3)}$ .

ANSWER.

$$(B.5.22) \quad \mathbf{K}^{(2,3)} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

□

**B.5.4. Kronecker Product and Vectorization in Tile Notation.** The Kronecker product of  $m \text{---} \mathbf{A} \text{---} n$  and  $r \text{---} \mathbf{B} \text{---} q$  is the following concatenation of  $\mathbf{A}$  and  $\mathbf{B}$  with members of a certain family of three-way arrays  $\Pi^{(i,j)}$ :

$$(B.5.23) \quad mr \text{---} \mathbf{A} \otimes \mathbf{B} \text{---} nq = mr \text{---} \Pi \begin{matrix} \text{---} m \\ \text{---} r \end{matrix} \begin{matrix} \mathbf{A} \\ \mathbf{B} \end{matrix} \text{---} nq$$

Strictly speaking we should have written  $\Pi^{(m,r)}$  and  $\Pi^{(n,q)}$  for the two  $\Pi$ -arrays (B.5.23), but the superscripts can be inferred from the context: the first superscript is the dimension of the Northeast arm, and the second that of the Southeast arm.

Vectorization uses a member of the same family  $\Pi^{(m,n)}$  to convert the matrix  $n \text{---} \mathbf{A} \text{---} m$  into the vector

$$(B.5.24) \quad mn \text{---} \operatorname{vec} \mathbf{A} = mn \text{---} \Pi \begin{matrix} \text{---} m \\ \text{---} n \end{matrix} \mathbf{A}$$

This equation is a little awkward because the  $\mathbf{A}$  is here a  $n \times m$  matrix, while elsewhere it is a  $m \times n$  matrix. It would have been more consistent with the lexicographical ordering used in the Kronecker product to define vectorization as the stacking of the row vectors; then some of the formulas would have looked more natural.

The array  $\Pi^{(m,n)} = mn \text{---} \Pi \begin{matrix} \text{---} m \\ \text{---} n \end{matrix}$  exists for every  $m \geq 1$  and  $n \geq 1$ . The

dimension of the West arm is always the product of the dimensions of the two East arms. The elements of  $\Pi^{(m,n)}$  will be given in (B.5.30) below; but first I will list three important properties of these arrays and give examples of their application.

First of all, each  $\Pi^{(m,n)}$  satisfies

$$(B.5.25) \quad \begin{matrix} m \\ \diagdown \\ \Pi \\ \diagup \\ n \end{matrix} \text{---} mn \text{---} \begin{matrix} m \\ \diagdown \\ \Pi \\ \diagup \\ n \end{matrix} = \begin{matrix} m & m \\ \diagdown & \diagup \\ & \\ \diagup & \diagdown \\ n & n \end{matrix}$$

Let us discuss the meaning of (B.5.25) in detail. The lefthand side of (B.5.25) shows the concatenation of two copies of the three-way array  $\Pi^{(m,n)}$  in a certain way that yields a 4-way array. Now look at the righthand side. The arm  $m \text{---} m$  by itself (which was bent only in order to remove any doubt about which arm to the left of the equal sign corresponds to which arm to the right) represents the neutral element under concatenation (i.e., the  $m \times m$  identity matrix). Writing two arrays next to each other without joining any arms represents their outer product, i.e., the array whose rank is the sum of the ranks of the arrays involved, and whose elements are all possible products of elements of the first array with elements of the second array.

The second identity satisfied by  $\Pi^{(m,n)}$  is

$$(B.5.26) \quad mn \text{---} \Pi \begin{array}{c} \text{---} m \\ \text{---} n \end{array} \Pi \text{---} mn = mn \text{---} mn .$$

Finally, there is also associativity:

$$(B.5.27) \quad mnp \text{---} \Pi \begin{array}{c} \text{---} m \\ \text{---} n \\ \text{---} p \end{array} \Pi \text{---} mnp = mnp \text{---} \Pi \begin{array}{c} \text{---} m \\ \text{---} n \\ \text{---} p \end{array} \Pi \text{---} mnp$$

Here is the answer to Problem 407 in tile notation:

$$(B.5.28) \quad \boxed{\text{tr } B C} = \boxed{B} \boxed{C} = \boxed{B} \Pi \Pi \boxed{C} = \boxed{\text{vec } B} \text{---} \boxed{\text{vec } C} = \boxed{(\text{vec } B)^T \text{vec } C}$$

Equation (B.5.25) was central for obtaining the result. The answer to Problem 410 also relies on equation (B.5.25):

$$(B.5.29) \quad \boxed{C^T \otimes A} \text{---} \boxed{\text{vec } B} = \Pi \begin{array}{c} \boxed{C} \\ \boxed{A} \end{array} \Pi \text{---} \Pi \boxed{B} = \Pi \begin{array}{c} \boxed{C} \\ \boxed{A} \end{array} \text{---} \boxed{B} = \boxed{\text{vec } ABC}$$

**B.5.5. Looking Inside the Kronecker Arrays.** It is necessary to open up the arrays from the  $\Pi$ -family and look at them “element by element,” in order to verify (B.5.23), (B.5.24), (B.5.25), (B.5.26), and (B.5.27). The elements of  $\Pi^{(m,n)}$ , which can be written in tile notation by saturating the array with unit vectors, are

$$(B.5.30) \quad \pi_{\theta\mu\nu}^{(m,n)} = \boxed{\theta} \text{---} mn \text{---} \Pi \begin{array}{c} \text{---} m \\ \text{---} n \end{array} \begin{array}{c} \boxed{\mu} \\ \boxed{\nu} \end{array} = \begin{cases} 1 & \text{if } \theta = (\mu - 1)n + \nu \\ 0 & \text{otherwise.} \end{cases}$$

Note that for every  $\theta$  there is exactly one  $\mu$  and one  $\nu$  such that  $\pi_{\theta\mu\nu}^{(m,n)} = 1$ ; for other values of  $\mu$  and  $\nu$ ,  $\pi_{\theta\mu\nu}^{(m,n)} = 0$ .

Writing  $\boxed{\nu} \text{---} \boxed{A} \text{---} \boxed{\mu} = a_{\nu\mu}$  and  $\boxed{\theta} \text{---} \boxed{\text{vec } A} = c_\theta$ , (B.5.24) reads

$$(B.5.31) \quad c_\theta = \sum_{\mu,\nu} \pi_{\theta\mu\nu}^{(m,n)} a_{\nu\mu},$$

which coincides with definition (B.5.2) of  $\text{vec } A$ .

One also checks that (B.5.23) is (B.5.3). Calling  $A \otimes B = C$ , it follows from (B.5.23) that

$$(B.5.32) \quad c_{\phi\theta} = \sum_{\mu,\nu,\rho,\kappa} \pi_{\phi\mu\rho}^{(m,r)} a_{\mu\nu} b_{\rho\kappa} \pi_{\theta\nu\kappa}^{(n,q)}.$$

For  $1 \leq \phi \leq r$  one gets a nonzero  $\pi_{\phi\mu\rho}^{(m,r)}$  only for  $\mu = 1$  and  $\rho = \phi$ , and for  $1 \leq \theta \leq q$  one gets a nonzero  $\pi_{\theta\nu\kappa}^{(n,q)}$  only for  $\nu = 1$  and  $\kappa = \theta$ . Therefore  $c_{\phi\theta} = a_{11} b_{\phi\theta}$  for elements of matrix  $C$  with  $\phi \leq r$  and  $\theta \leq q$ . Etc.

The proof of (B.5.25) uses the fact that for every  $\theta$  there is exactly one  $\mu$  and one  $\nu$  such that  $\pi_{\theta\mu\nu}^{(m,n)} \neq 0$ :

$$(B.5.33) \quad \sum_{\theta=1}^{\theta=mn} \pi_{\theta\mu\nu}^{(m,n)} \pi_{\theta\omega\sigma}^{(m,n)} = \begin{cases} 1 & \text{if } \mu = \omega \text{ and } \nu = \sigma \\ 0 & \text{otherwise} \end{cases}$$

Similarly, (B.5.26) and (B.5.27) can be shown by elementary but tedious proof. The best verification of these rules is their implementation in a computer language; see Section ?? below.

**B.5.6. The Commutation Matrix in Tile Notation.** The simplest way to represent the commutation matrix  $K^{(m,n)}$  in a tile is

$$(B.5.34) \quad K^{(m,n)} = mn \text{---} \Pi \begin{array}{c} \text{---} m \\ \text{---} n \end{array} \Pi \text{---} mn .$$

This should not be confused with the lefthand side of (B.5.26):  $K^{(m,n)}$  is composed of  $\Pi^{(m,n)}$  on its West and  $\Pi^{(n,m)}$  on its East side, while (B.5.26) contains  $\Pi^{(m,n)}$  twice. We will therefore use the following representation, mathematically equivalent to (B.5.34), which makes it easier to see the effects of  $K^{(m,n)}$ :

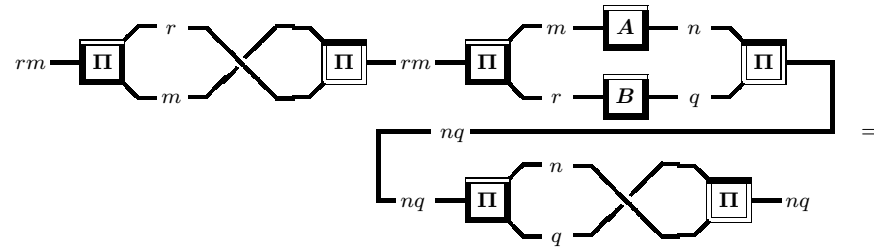
$$(B.5.35) \quad K^{(m,n)} = mn \text{---} \Pi \begin{array}{c} \text{---} m \\ \text{---} n \end{array} \Pi \text{---} mn .$$

PROBLEM 415. Using the definition (B.5.35) show that  $\mathbf{K}^{(m,n)}\mathbf{K}^{(n,m)} = \mathbf{I}_{mn}$ , the  $mn \times mn$  identity matrix.

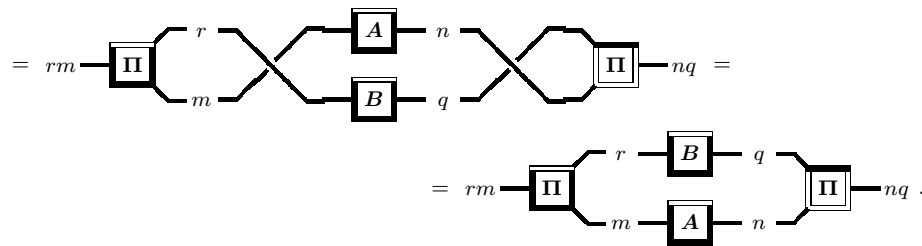
ANSWER. You will need (B.5.25) and (B.5.26). □

PROBLEM 416. Prove (B.5.21) in tile notation.

ANSWER. Start with a tile representation of  $\mathbf{K}^{(r,m)}(\mathbf{A} \otimes \mathbf{B})\mathbf{K}^{(n,q)}$ :



Now use (B.5.25) twice to get



□



## APPENDIX C

## Matrix Differentiation

## C.1. First Derivatives

Let us first consider the scalar case and then generalize from there. The derivative of a function  $f$  is often written

$$(C.1.1) \quad \frac{dy}{dx} = f'(x)$$

Multiply through by  $dx$  to get  $dy = f'(x) dx$ . In order to see the meaning of this equation, we must know the definition  $dy = f(x + dx) - f(x)$ . Therefore one obtains  $f(x + dx) = f(x) + f'(x) dx$ . If one holds  $x$  constant and only varies  $dx$  this formula shows that in an infinitesimal neighborhood of  $x$ , the function  $f$  is an *affine* function of  $dx$ , i.e., a linear function of  $dx$  with a constant term:  $f(x)$  is the intercept, i.e., the value for  $dx = 0$ , and  $f'(x)$  is the slope parameter.

Now let us transfer this argument to vector functions  $\mathbf{y} = \mathbf{f}(\mathbf{x})$ . Here  $\mathbf{y}$  is a  $n$ -vector and  $\mathbf{x}$  a  $m$ -vector, i.e.,  $\mathbf{f}$  is a  $n$ -tuple of functions of  $m$  variables each

$$(C.1.2) \quad \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_m) \\ \vdots \\ f_n(x_1, \dots, x_m) \end{bmatrix}$$

One may also say,  $\mathbf{f}$  is a  $n$ -vector, each element of which depends on  $\mathbf{x}$ . Again, under certain differentiability conditions, it is possible to write this function infinitesimally as an affine function, i.e., one can write

$$(C.1.3) \quad \mathbf{f}(\mathbf{x} + d\mathbf{x}) = \mathbf{f}(\mathbf{x}) + \mathbf{A}d\mathbf{x}.$$

Here the coefficient of  $d\mathbf{x}$  is no longer a scalar but necessarily a *matrix*  $\mathbf{A}$  (whose elements again depend on  $\mathbf{x}$ ).  $\mathbf{A}$  is called the *Jacobian matrix* of  $\mathbf{f}$ . The Jacobian matrix generalizes the concept of a derivative to vectors. Instead of a prime denoting the derivative, as in  $f'(x)$ , one writes  $\mathbf{A} = \mathbf{D}\mathbf{f}$ .

PROBLEM 417. 2 points If  $f$  is a scalar function of a vector argument  $\mathbf{x}$ , is its Jacobian matrix  $\mathbf{A}$  a row vector or a column vector? Explain why this must be so.

The Jacobian  $\mathbf{A}$  defined in this way turns out to have a very simple functional form: its elements are the partial derivatives of all components of  $\mathbf{f}$  with respect to all components of  $\mathbf{x}$ :

$$(C.1.4) \quad a_{ij} = \frac{\partial f_i}{\partial x_j}.$$

Since in this matrix  $\mathbf{f}$  acts as column and  $\mathbf{x}$  as a row vector, this matrix can be written, using matrix differentiation notation, as  $\mathbf{A}(\mathbf{x}) = \partial \mathbf{f}(\mathbf{x}) / \partial \mathbf{x}^\top$ .

Strictly speaking, *matrix* notation can be used for matrix differentiation only if we differentiate a column vector (or scalar) with respect to a row vector (or scalar) or if we differentiate a scalar with respect to a matrix or a matrix with respect to a scalar. If we want to differentiate matrices with respect to vectors or vectors with respect to matrices or matrices with respect to each other, we need the tile notation for arrays. A different, much less enlightening approach is to first “vectorize” the matrices involved. Both of those methods will be discussed later.

If the dependence of  $\mathbf{y}$  on  $\mathbf{x}$  can be expressed in terms of matrix operations or more general array concatenations, then some useful matrix differentiation rules exist.

The simplest matrix differentiation rule, for  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  with

$$(C.1.5) \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \quad \text{and} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

is

$$(C.1.6) \quad \partial \mathbf{w}^\top \mathbf{x} / \partial \mathbf{x}^\top = \mathbf{w}^\top$$

Here is the proof of (C.1.6):

$$\begin{aligned} \frac{\partial \mathbf{w}^\top \mathbf{x}}{\partial \mathbf{x}^\top} &= \left[ \frac{\partial}{\partial x_1} (w_1 x_1 + \dots + w_n x_n) \quad \dots \quad \frac{\partial}{\partial x_n} (w_1 x_1 + \dots + w_n x_n) \right] \\ &= [w_1 \quad \dots \quad w_n] = \mathbf{w}^\top \end{aligned}$$

The second rule, for  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{M} \mathbf{x}$  and  $\mathbf{M}$  symmetric, is:

$$(C.1.7) \quad \partial \mathbf{x}^\top \mathbf{M} \mathbf{x} / \partial \mathbf{x}^\top = 2\mathbf{x}^\top \mathbf{M}.$$



Here is a differentiation rule for a matrix with respect to a matrix, first written element by element, and then in tiles: If  $Y = AXB$ , i.e.,  $y_{im} = \sum_{j,k} a_{ij}x_{jk}b_{km}$ , then  $\frac{\partial y_{im}}{\partial x_{jk}} = a_{ij}a_{km}$ , because for every fixed  $i$  and  $m$  this sum contains only one term which has  $x_{jk}$  in it, namely,  $a_{ij}x_{jk}b_{km}$ . In tiles:

$$(C.1.18) \quad \begin{array}{c} \text{---} \\ \boxed{A} \\ \text{---} \\ \boxed{X} \\ \text{---} \\ \boxed{B} \\ \text{---} \end{array} / \frac{\partial}{\partial} \begin{array}{c} \text{---} \\ \boxed{X} \\ \text{---} \end{array} = \begin{array}{c} \boxed{A} \\ \text{---} \\ \boxed{B} \\ \text{---} \end{array}$$

Equations (C.1.17) and (C.1.18) can be obtained from (C.1.12) and (C.1.15) by extended substitution, since a bundle of several arms can always be considered as one arm. For instance, (C.1.17) can be written

$$\frac{\partial}{\partial} \begin{array}{c} \text{---} \\ \boxed{A} \\ \text{---} \\ \text{---} \\ \text{---} \\ \boxed{X} \\ \text{---} \\ \text{---} \end{array} / \frac{\partial}{\partial} \begin{array}{c} \text{---} \\ \boxed{X} \\ \text{---} \end{array} = \begin{array}{c} \boxed{A} \\ \text{---} \\ \text{---} \end{array}$$

and this is a special case of (C.1.12), since the two parallel arms can be treated as one arm. With a better development of the logic underlying this notation, it will not be necessary to formulate them as separate theorems; all matrix differentiation rules given so far are trivial applications of (C.1.15).

PROBLEM 419. As a special case of (C.1.18) show that  $\frac{\partial x^\top A y}{\partial A^\top} = yx^\top$ .

ANSWER.

$$(C.1.19) \quad \begin{array}{c} \boxed{x} \\ \text{---} \\ \boxed{A} \\ \text{---} \\ \boxed{y} \\ \text{---} \end{array} / \frac{\partial}{\partial} \begin{array}{c} \text{---} \\ \boxed{A} \\ \text{---} \end{array} = \begin{array}{c} \boxed{x} \\ \text{---} \\ \boxed{y} \\ \text{---} \end{array}$$

Here is a basic differentiation rule for *bilinear* array concatenations: if

$$(C.1.20) \quad \text{---} \boxed{y} = \text{---} \boxed{A} \begin{array}{c} \text{---} \boxed{x} \\ \text{---} \boxed{x} \end{array}$$

then one gets the following simple generalization of (C.1.13):

$$(C.1.21) \quad \frac{\partial}{\partial} \begin{array}{c} \text{---} \boxed{A} \\ \text{---} \boxed{x} \\ \text{---} \boxed{x} \end{array} / \frac{\partial}{\partial} \begin{array}{c} \text{---} \boxed{x} \\ \text{---} \end{array} = \begin{array}{c} \text{---} \boxed{A} \\ \text{---} \boxed{x} \end{array} + \begin{array}{c} \text{---} \boxed{A} \\ \text{---} \boxed{x} \end{array}$$

PROOF.  $y_i = \sum_{j,k} a_{ijk}x_jx_k$ . For a given  $i$ , this has  $x_p^2$  in the term  $a_{ipp}x_p^2$ , and it has  $x_p$  in the terms  $a_{ipk}x_px_k$  where  $p \neq k$ , and in  $a_{ijp}x_jx_p$  where  $j \neq p$ . The derivatives of these terms are  $2a_{ipp}x_p + \sum_{k \neq p} a_{ipk}x_k + \sum_{j \neq p} a_{ijp}x_j$ , which simplifies to  $\sum_k a_{ipk}x_k + \sum_j a_{ijp}x_j$ . This is the  $i, p$ -element of the matrix on the rhs of (C.1.21).

But there are also other ways to have the array  $X$  occur twice in a concatenation  $Y$ . If  $Y = X^\top X$  then  $y_{ik} = \sum_j x_{ji}x_{jk}$  and therefore  $\partial y_{ik}/\partial x_{lm} = 0$  if  $m \neq i$  and  $m \neq k$ . Now assume  $m = i \neq k$ :  $\partial y_{ik}/\partial x_{li} = \partial x_{li}x_{lk}/\partial x_{li} = x_{lk}$ . Now assume  $m = k \neq i$ :  $\partial y_{ik}/\partial x_{lk} = \partial x_{li}x_{lk}/\partial x_{lk} = x_{li}$ . And if  $m = k = i$  then one gets the sum of the two above:  $\partial y_{ii}/\partial x_{li} = \partial x_{li}^2/\partial x_{li} = 2x_{li}$ . In tiles this is

$$(C.1.22) \quad \frac{\partial X^\top X}{\partial X^\top} = \frac{\partial}{\partial} \begin{array}{c} \boxed{X} \\ \text{---} \\ \boxed{X} \\ \text{---} \end{array} / \frac{\partial}{\partial} \begin{array}{c} \text{---} \boxed{X} \\ \text{---} \end{array} = \begin{array}{c} \text{---} \boxed{X} \\ \text{---} \end{array} + \begin{array}{c} \text{---} \boxed{X} \\ \text{---} \end{array}$$

This rule is helpful for differentiating the multivariate Normal likelihood function

A computer implementation of this tile notation should contain algorithms that automatically take the derivatives of these array concatenations.

Here are some more matrix differentiation rules:

Chain rule: If  $g = g(\eta)$  and  $\eta = \eta(\beta)$  are two vector functions, then

$$(C.1.23) \quad \partial g / \partial \beta^\top = \partial g / \partial \eta^\top \cdot \partial \eta / \partial \beta^\top$$

For instance, the linear least squares objective function is  $SSE = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}$  where  $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\beta$ . Application of the chain rule gives  $\partial SSE / \partial \beta^\top = \partial SSE / \partial \hat{\boldsymbol{\varepsilon}}^\top \cdot \partial \hat{\boldsymbol{\varepsilon}} / \partial \beta^\top = 2\hat{\boldsymbol{\varepsilon}}^\top (-\mathbf{X})$  which is the same result as in (14.2.2).

If  $A$  is nonsingular then

$$(C.1.24) \quad \frac{\partial \log \det A}{\partial A^\top} = A^{-1}$$

Proof in [Gre97, pp. 52/3].

## Bibliography

- [AD75] J. Aczél and Z. Daróczy. *On Measures of Information and their Characterizations*. Academic Press, 1975.
- [Alb69] Arthur E. Albert. Conditions for positive and negative semidefiniteness in terms of pseudoinverses. *SIAM (Society for Industrial and Applied Mathematics) Journal of Applied Mathematics*, 17:434–440, 1969.
- [Alb72] Arthur E. Albert. *Regression and the Moore-Penrose Pseudoinverse*. Academic Press, New York and London, 1972.
- [Ame85] Takeshi Amemiya. *Advanced Econometrics*. Harvard University Press, 1985.
- [Ame94] Takeshi Amemiya. *Introduction to Statistics and Econometrics*. Harvard University Press, Cambridge, MA, 1994.
- [Bar82] Vic Barnett. *Comparative Statistical Inference*. Wiley, New York, 1982.
- [BCW96] Richard A. Becker, John M. Chambers, and Allan R. Wilks. *The New S Language: A Programming Environment for Data Analysis and Graphics*. Chapman and Hall, 1996. Reprint of the 1988 Wadsworth edition.
- [BD77] Peter J. Bickel and Kjell A. Doksum. *Mathematical Statistics: Basic Ideas and Selected Topics*. Holden-Day, San Francisco, 1977.
- [Ber91] Ernst R. Berndt. *The Practice of Econometrics: Classic and Contemporary*. Addison-Wesley, Reading, Massachusetts, 1991.
- [BKW80] David A. Belsley, Edwin Kuh, and Roy E. Welsch. *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. Wiley, New York, 1980.
- [CD28] Charles W. Cobb and Paul H. Douglas. A theory of production. *American Economic Review*, 18(1, Suppl.):139–165, 1928. J.
- [CD97] Wojciech W. Charemza and Derek F. Deadman. *New Directions in Econometric Practice: General to Specific Modelling, Cointegration, and Vector Autoregression*. Edward Elgar, Cheltenham, UK; Lynne, NH, 2nd ed. edition, 1997.
- [CH93] John M. Chambers and Trevor Hastie, editors. *Statistical Models in S*. Chapman and Hall, 1993.
- [Cho60] G. C. Chow. Tests of equality between sets of coefficients in two linear regressions. *Econometrica*, 28:591–605, July 1960.
- [Chr87] Ronald Christensen. *Plane Answers to Complex Questions; The Theory of Linear Models*. Springer-Verlag, New York, 1987.
- [Coh50] A. C. Cohen. Estimating the mean and variance of normal populations from singly and doubly truncated samples. *Annals of Mathematical Statistics*, pages 557–569, 1950.
- [Coo77] R. Dennis Cook. Detection of influential observations in linear regression. *Technometrics*, 19(1):15–18, February 1977.
- [Coo98] R. Dennis Cook. *Regression Graphics: Ideas for Studying Regressions through Graphics*. Series in Probability and Statistics. Wiley, New York, 1998.

- [Cor69] J. Cornfield. The Bayesian outlook and its applications. *Biometrics*, 25:617–657, 1969.
- [Cow77] Frank Alan Cowell. *Measuring Inequality: Techniques for the Social Sciences*. Wiley, New York, 1977.
- [Cra43] A. T. Craig. Note on the independence of certain quadratic forms. *Annals of Mathematical Statistics*, 14:195, 1943.
- [CT91] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Series in Telecommunications. Wiley, New York, 1991.
- [CW99] R. Dennis Cook and Sanford Weisberg. *Applied Regression Including Computing Graphics*. Wiley, 1999.
- [Daw79a] A. P. Dawid. Conditional independence in statistical theory. *JRSS(B)*, 41(1):1–31, 1979.
- [Daw79b] A. P. Dawid. Some misleading arguments involving conditional independence. *JRSS(B)*, 41(2):249–252, 1979.
- [Daw80] A. P. Dawid. Conditional independence for statistical operations. *Annals of Statistics*, 8:598–617, 1980.
- [Dhr86] Phoebus J. Dhrymes. Limited dependent variables. In Zvi Griliches and Michael Intriligator, editors, *Handbook of Econometrics*, volume 3, chapter 27, pages 1567–1600. North-Holland, Amsterdam, 1986.
- [DL91] Gerard Dumenil and Dominique Levy. The U.S. economy since the Civil War: Sources and construction of the series. Technical report, CEPREMAP, LAREA-CEDRA, November 1991.
- [DM93] Russell Davidson and James G. MacKinnon. *Estimation and Inference in Econometrics*. Oxford University Press, New York, 1993.
- [Dou92] Christopher Dougherty. *Introduction to Econometrics*. Oxford University Press, Oxford, 1992.
- [DP20] R. E. Day and W. M. Persons. An index of the physical volume of production. *Review of Economic Statistics*, II:309–37, 361–67, 1920.
- [Fis] R. A. Fisher. Theory of statistical estimation. *Proceedings of the Cambridge Philosophical Society*, 22.
- [Fri57] Milton Friedman. *A Theory of the Consumption Function*. Princeton University Press, Princeton, 1957.
- [FS91] Milton Friedman and Anna J. Schwarz. Alternative approaches to analyzing economic data. *American Economic Review*, 81(1):39–49, March 1991.
- [GJM96] Amos Golan, George Judge, and Douglas Miller. *Maximum Entropy Econometrics: Robust Estimation with Limited Data*. Wiley, Chichester, England, 1996.
- [Gra76] Franklin A. Graybill. *Theory and Application of the Linear Model*. Duxbury Press, North Scituate, Mass., 1976.
- [Gra83] Franklin A. Graybill. *Matrices with Applications in Statistics*. Wadsworth and Brooks/Cole, Pacific Grove, CA, second edition, 1983.
- [Gre97] William H. Greene. *Econometric Analysis*. Prentice Hall, Upper Saddle River, NJ, third edition, 1997.
- [Gre03] William H. Greene. *Econometric Analysis*. Prentice Hall, Upper Saddle River, NJ, Jersey 07458, fifth edition, 2003.
- [Hal78] Robert E. Hall. Stochastic implications of the life cycle-permanent income hypothesis: Theory and evidence. *Journal of Political Economy*, pages 971–987, December 1978.
- [Hen95] David F. Hendry. *Dynamic Econometrics*. Oxford University Press, Oxford, New York, 1995.
- [HK79] James M. Henle and Eugene M. Kleinberg. *Infinitesimal Calculus*. MIT Press, 1979.

- [Hou51] H. S. Houthakker. Some calculations on electricity consumption in Great Britain. *Journal of the Royal Statistical Society (A)*, (114 part III):351–371, 1951. J.
- [HT83] Robert V. Hogg and Elliot A. Tanis. *Probability and Statistical Inference*. Macmillan, second edition, 1983.
- [HVdP02] Ben J. Hejdra and Frederick Van der Ploeg. *Foundations of Modern Macroeconomics*. Oxford University Press, 2002.
- [JHG+88] George G. Judge, R. Carter Hill, William E. Griffiths, Helmut Lütkepohl, and Tsoung-Chao Lee. *Introduction to the Theory and Practice of Econometrics*. Wiley, New York, second edition, 1988.
- [JK70] Norman Johnson and Samuel Kotz. *Continuous Univariate Distributions*, volume 1. Houghton Mifflin, Boston, 1970.
- [KA69] J. Koerts and A. P. J. Abramanse. *On the Theory and Application of the General Linear Model*. Rotterdam University Press, Rotterdam, 1969.
- [Kap89] Jagat Narain Kapur. *Maximum Entropy Models in Science and Engineering*. Wiley, 1989.
- [Ken98] Peter Kennedy. *A Guide to Econometrics*. MIT Press, Cambridge, MA, fourth edition, 1998.
- [Khi57] R. T. Khinchin. *Mathematical Foundations of Information Theory*. Dover Publications, New York, 1957.
- [Kme86] Jan Kmenta. *Elements of Econometrics*. Macmillan, New York, second edition, 1986.
- [Knu81] Donald E. Knuth. *Seminumerical Algorithms*, volume 2 of *The Art of Computer Programming*. Addison-Wesley, second edition, 1981.
- [Krz88] W. J. Krzanowski. *Principles of Multivariate Analysis: A User's Perspective*. Clarendon Press, Oxford, 1988.
- [KS79] Sir Maurice Kendall and Alan Stuart. *The Advanced Theory of Statistics*, volume 2. Griffin, London, fourth edition, 1979.
- [Ksh19] Anant M. Kshirsagar. *Multivariate Analysis*. Marcel Dekker, New York and Basel, 19??
- [Lan69] H. O. Lancaster. *The Chi-Squared Distribution*. Wiley, 1969.
- [Lar82] Harold Larson. *Introduction to Probability and Statistical Inference*. Wiley, 1982.
- [Lea75] Edward E. Leamer. A result on the sign of the restricted least squares estimator. *Journal of Econometrics*, 3:387–390, 1975.
- [Mag88] Jan R. Magnus. *Linear Structures*. Oxford University Press, New York, 1988.
- [Mal80] E. Malinvaud. *Statistical Methods of Econometrics*. North-Holland, Amsterdam, third edition, 1980.
- [MN88] Jan R. Magnus and Heinz Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Wiley, Chichester, 1988.
- [Mor65] A. Q. Morton. The authorship of Greek prose (with discussion). *Journal of the Royal Statistical Society, Series A*, 128:169–233, 1965.
- [Mor73] Trenchard More, Jr. Axioms and theorems for a theory of arrays. *IBM Journal of Research and Development*, 17(2):135–175, March 1973.
- [Mor02] Jamie Morgan. The global power of orthodox economics. *Journal of Critical Realism*, 1(2):7–34, May 2002.
- [MR91] Ieke Moerdijk and Gonzalo E. Reyes. *Models for Smooth Infinitesimal Analysis*. Springer-Verlag, New York, 1991.
- [MS86] Parry Hiram Moon and Domina Eberle Spencer. *Theory of Holors; A Generalization of Tensors*. Cambridge University Press, 1986.

- [Rao62] C. Radhakrishna Rao. A note on a generalized inverse of a matrix with applications to problems in mathematical statistics. *Journal of the Royal Statistical Society, Series B*, 24:152–158, 1962.
- [Rao73] C. Radhakrishna Rao. *Linear Statistical Inference and Its Applications*. Wiley, New York, second edition, 1973.
- [Rao97] C. Radhakrishna Rao. *Statistics and Truth: Putting Chance to Work*. World Scientific, Singapore, second edition, 1997.
- [Rei89] Rolf-Dieter Reiss. *Approximate Distributions of Order Statistics*. Springer-Verlag, New York, 1989.
- [Rén70] Alfred Rényi. *Foundations of Probability*. Holden-Day, San Francisco, 1970.
- [Rie77] E. Rietsch. The maximum entropy approach to inverse problems. *Journal of Geophysical Research*, 82:42, 1977.
- [Rie85] E. Rietsch. On an alleged breakdown of the maximum-entropy principle. In C. Ray Smith and Jr W. T. Grandy, editors, *Maximum-Entropy and Bayesian Methods in Inverse Problems*, pages 67–82. D. Reidel, Dordrecht, Boston, Lancaster, 1985.
- [Rob70] Herbert Robbins. Statistical methods related to the law of the iterated logarithm. *Annals of Mathematical Statistics*, 41:1397–1409, 1970.
- [Rob74] Abraham Robinson. *Non-Standard Analysis*. North Holland, Amsterdam, 1974.
- [Ron02] Amit Ron. Regression analysis and the philosophy of social science: A critical realist view. *Journal of Critical Realism*, 1(1):119–142, November 2002.
- [Roy97] Richard M. Royall. *Statistical evidence: A Likelihood Paradigm*. Number 71 in *Monographs on Statistics and Applied Probability*. Chapman & Hall, London; New York, 1997.
- [RZ78] L. S. Robertson and P. L. Zador. Driver education and fatal crash involvement of teen drivers. *American Journal of Public Health*, 68:959–65, 1978.
- [Seb77] G. A. F. Seber. *Linear Regression Analysis*. Wiley, New York, 1977.
- [Sel58] H. C. Selvin. Durkheim's suicide and problems of empirical research. *American Journal of Sociology*, 63:607–619, 1958.
- [SG85] John Skilling and S. F. Gull. Algorithms and applications. In C. Ray Smith and Jr W. T. Grandy, editors, *Maximum-Entropy and Bayesian Methods in Inverse Problems*, pages 83–132. D. Reidel, Dordrecht, Boston, Lancaster, 1985.
- [SM86] Hans Schneeweiß and Hans-Joachim Mittag. *Lineare Modelle mit fehlerbehafteten Daten*. Physica Verlag, Heidelberg, Wien, 1986.
- [Spr98] Peter Sprent. *Data Driven Statistical Methods*. Texts in statistical science. Chapman & Hall, London; New York, 1st ed. edition, 1998.
- [SS35] J. A. Schouten and Dirk J. Struik. *Einführung in die neuen Methoden der Differentialgeometrie*, volume I. 1935.
- [SW76] Thomas J. Sargent and Neil Wallace. Rational expectations and the theory of economic policy. *Journal of Monetary Economics*, 2:169–183, 1976.
- [The71] Henri Theil. *Principles of Econometrics*. Wiley, New York, 1971.
- [Tin51] J. Tinbergen. *Econometrics*. George Allen & Unwin Ltd., London, 1951.
- [TS61] Henri Theil and A. Schweitzer. The best quadratic estimator of the residual variance in regression analysis. *Statistica Neerlandica*, 15:19–23, 1961.
- [Wit85] Uli Wittman. *Das Konzept rationaler Preiserwartungen*, volume 241 of *Lecture Notes in Economics and Mathematical Systems*. Springer, 1985.
- [Yul07] G. V. Yule. On the theory of correlation for any number of variables treated by a new system of notation. *Proc. Roy. Soc. London A*, 79:182, 1907.