

SECOND MIDTERM ECON 7800 FALL 2003

ECONOMICS DEPARTMENT, UNIVERSITY OF UTAH

Problem 21. *2 points The conditional density is the joint divided by the marginal:*

$$(1) \quad f_{y|x}(y, x) = \frac{f_{x,y}(x, y)}{f_x(x)}.$$

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,

$$(2) \quad \int_{y=-\infty}^{+\infty} f_{y|x=x}(y, x) dy = \frac{\int_{y=-\infty}^{+\infty} f_{x,y}(x, y) dy}{f_x(x)} = \frac{f_x(x)}{f_x(x)} = 1$$

because of the formula for the marginal:

$$(3) \quad f_x(x) = \int_{y=-\infty}^{+\infty} f_{x,y}(x,y) dy$$

You see that formula (1) divides the joint density exactly by the right number which makes the integral equal to 1. \square

Problem 22. Assume the vector $\mathbf{x} = [x_1, \dots, x_j]^\top$ and the scalar y are jointly distributed random variables, and assume conditional means exist. Define $\varepsilon = y - \mathbf{E}[y|\mathbf{x}]$.

For this problem you may use (1) the theorem of iterated expectations $\mathbf{E}[\mathbf{E}[y|\mathbf{x}]] = \mathbf{E}[y]$, (2) the additivity $\mathbf{E}[g(y) + h(y)|\mathbf{x}] = \mathbf{E}[g(y)|\mathbf{x}] + \mathbf{E}[h(y)|\mathbf{x}]$, and (3) the fact that $\mathbf{E}[g(\mathbf{x})h(y)|\mathbf{x}] = g(\mathbf{x})\mathbf{E}[h(y)|\mathbf{x}]$.

• **a.** 5 points Demonstrate the following identities:

$$(4) \quad \mathbf{E}[\varepsilon|\mathbf{x}] = 0$$

$$(5) \quad \mathbf{E}[\varepsilon] = 0$$

$$(6) \quad \mathbf{E}[x_i \varepsilon|\mathbf{x}] = 0 \quad \text{for all } i, 1 \leq i \leq j$$

$$(7) \quad \mathbf{E}[x_i \varepsilon] = 0 \quad \text{for all } i, 1 \leq i \leq j$$

$$(8) \quad \text{cov}[\mathbf{x}_i, \varepsilon] = 0 \quad \text{for all } i, 1 \leq i \leq j.$$

Interpretation of (8): ε is the error in the best prediction of y based on \mathbf{x} . If this error were correlated with one of the components x_i , then this correlation could be used to construct a better prediction of y .

Answer. (4): $E[\varepsilon|\mathbf{x}] = E[y|\mathbf{x}] - E[E[y|\mathbf{x}]|\mathbf{x}] = 0$ since $E[y|\mathbf{x}]$ 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 but the law that the expected value of a constant is a constant.)

(5) follows from (4) (i.e., (4) is stronger than (5)): 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 (4): $E[\varepsilon] = E[y - E[y|\mathbf{x}]] = 0$. Either way you need the law of iterated expectations for this.

$$(6): E[x_i \varepsilon|\mathbf{x}] = x_i E[\varepsilon|\mathbf{x}] = 0.$$

(7): $E[x_i \varepsilon] = E[E[x_i \varepsilon|\mathbf{x}]] = E[0] = 0$; or in one swoop: $E[x_i \varepsilon] = E[x_i y - x_i E[y|\mathbf{x}]] = E[x_i y - E[x_i y|\mathbf{x}]] = E[x_i y] - E[x_i y] = 0$. The following “proof” is not correct: $E[x_i \varepsilon] = E[x_i] E[\varepsilon] = E[x_i] \cdot 0 = 0$. x_i and ε 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.

$$(8): \text{cov}[x_i, \varepsilon] = E[x_i \varepsilon] - E[x_i] E[\varepsilon] = 0 - E[x_i] \cdot 0 = 0. \quad \square$$

• **b.** 2 points If \mathbf{x} and y are jointly normal, show that \mathbf{x} and ε are independent, and that the variance of ε 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 ε 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 23. 4 points Let y_1, \dots, y_n be an arbitrary vector and α an arbitrary number. As usual, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. Show that

$$(9) \quad \sum_{i=1}^n (y_i - \alpha)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \alpha)^2$$

Answer.

$$(10) \quad \sum_{i=1}^n (y_i - \alpha)^2 = \sum_{i=1}^n ((y_i - \bar{y}) + (\bar{y} - \alpha))^2$$

$$(11) \quad = \sum_{i=1}^n (y_i - \bar{y})^2 + 2 \sum_{i=1}^n ((y_i - \bar{y})(\bar{y} - \alpha)) + \sum_{i=1}^n (\bar{y} - \alpha)^2$$

$$(12) \quad = \sum_{i=1}^n (y_i - \bar{y})^2 + 2(\bar{y} - \alpha) \sum_{i=1}^n (y_i - \bar{y}) + n(\bar{y} - \alpha)^2$$

Since the middle term is zero, (9) follows. \square

Problem 24. 4 points Show that

$$(13) \quad s_u^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

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) with $\alpha = E[y]$:

$$(14) \quad E\left[\sum_{i=1}^n (y_i - \bar{y})^2\right] = \sum_{i=1}^n E[(y_i - \mu)^2] - n E[(\bar{y} - \mu)^2]$$

$$(15) \quad = \sum_{i=1}^n \sigma^2 - n \frac{\sigma^2}{n} = (n-1)\sigma^2.$$

You do not need Normality for this. □

Problem 25.

- **a.** 2 points Verify that the matrix $D = I - \frac{1}{n}\mathbf{u}\mathbf{u}^\top$ is symmetric and idempotent.
- **b.** 1 point Compute the trace $\text{tr } D$.

Answer. $\text{tr } \mathbf{D} = n - 1$. One can see this either by writing down the matrix element by element, or use the linearity of the trace plus the rule that $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$. $\text{tr } \mathbf{I} = n$ and $\text{tr}(\boldsymbol{\iota} \boldsymbol{\iota}^\top) = \text{tr}(\boldsymbol{\iota}^\top \boldsymbol{\iota}) = \text{tr } n = n$. \square

• **c.** 1 point For any vector of observations \mathbf{y} compute $\mathbf{D}\mathbf{y}$.

Answer. Element by element one can write

$$(16) \quad \mathbf{D}\mathbf{y} = \begin{bmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \cdots & -\frac{1}{n} \\ \vdots & & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & \cdots & 1 - \frac{1}{n} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{bmatrix}$$

There is also a more elegant matrix theoretical proof available \square

• **d.** 1 point Is there a vector $\mathbf{a} \neq \mathbf{o}$ for which $\mathbf{D}\mathbf{a} = \mathbf{o}$? If so, give an example of such a vector.

Answer. $\boldsymbol{\iota}$ is, up to a scalar factor, the only nonzero vector with $\mathbf{D}\boldsymbol{\iota} = \mathbf{o}$. \square

• **e.** 1 point Show that the sample variance of a vector of observations \mathbf{y} can be written in matrix notation as

$$(17) \quad \text{The sample variance of } \mathbf{y} \text{ is } \frac{1}{n} \sum (y_i - \bar{y})^2 = \frac{1}{n} \mathbf{y}^\top \mathbf{D}\mathbf{y}$$

Answer. Let's get rid of the factor $\frac{1}{n}$ which appears on both sides: we have to show that

$$(18) \quad \sum (y_i - \bar{y})^2 = \mathbf{y}^\top \mathbf{D} \mathbf{y} = \mathbf{y}^\top \mathbf{D}^\top \mathbf{D} \mathbf{y}$$

This is the squared length of the vector $\mathbf{D} \mathbf{y}$ which we computed in part c. □

Problem 26. [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 t be an estimator of the nonrandom scalar parameter θ . $E[t - \theta]$ is called the bias of t , and $E[(t - \theta)^2]$ is called the mean squared error of t as an estimator of θ , written $\text{MSE}[t; \theta]$. Show that the MSE is the variance plus the squared bias, i.e., that*

$$(19) \quad \text{MSE}[t; \theta] = \text{var}[t] + (E[t - \theta])^2.$$

Answer. The most elegant proof, which also indicates what to do when θ is random, is:

$$(20) \quad \text{MSE}[t; \theta] = E[(t - \theta)^2] = \text{var}[t - \theta] + (E[t - \theta])^2 = \text{var}[t] + (E[t - \theta])^2.$$

□

• **b.** 2 points *For the rest of this problem assume that t is an unbiased estimator of θ with $\text{var}[t] > 0$. We will investigate whether one can get a better MSE if one*

estimates θ by a constant multiple at instead of t . Show that

$$(21) \quad \text{MSE}[at; \theta] = a^2 \text{var}[t] + (a - 1)^2 \theta^2.$$

Answer. $\text{var}[at] = a^2 \text{var}[t]$ and the bias of at is $E[at - \theta] = (a - 1)\theta$. Now apply (20). □

• **c.** 1 point Show that, whenever $a > 1$, then $\text{MSE}[at; \theta] > \text{MSE}[t; \theta]$. If one wants to decrease the MSE, one should therefore not choose $a > 1$.

Answer. $\text{MSE}[at; \theta] - \text{MSE}[t; \theta] = (a^2 - 1) \text{var}[t] + (a - 1)^2 \theta^2 > 0$ since $a > 1$ and $\text{var}[t] > 0$. □

• **d.** 2 points Show that

$$(22) \quad \left. \frac{d}{da} \text{MSE}[at; \theta] \right|_{a=1} > 0.$$

From this follows that the MSE of at is smaller than the MSE of t , as long as $a < 1$ and close enough to 1.

Answer. The derivative of (21) is

$$(23) \quad \frac{d}{da} \text{MSE}[at; \theta] = 2a \text{var}[t] + 2(a - 1)\theta^2$$

Plug $a = 1$ into this to get $2 \text{var}[t] > 0$. □

- **e.** 2 points By solving the first order condition show that the factor a which gives smallest MSE is

$$(24) \quad a = \frac{\theta^2}{\text{var}[t] + \theta^2}.$$

Answer. Rewrite (23) as $2a(\text{var}[t] + \theta^2) - 2\theta^2$ and set it zero. □

- **f.** 1 point Assume t has an exponential distribution with parameter $\lambda > 0$, i.e.,

$$(25) \quad f_t(t) = \lambda \exp(-\lambda t), \quad t \geq 0 \quad \text{and} \quad f_t(t) = 0 \quad \text{otherwise.}$$

Check that $f_t(t)$ is indeed a density function.

Answer. Since $\lambda > 0$, $f_t(t) > 0$ 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 the exponential distribution) prove that t is an unbiased estimator of $1/\lambda$, with $\text{var}[t] = 1/\lambda^2$.

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 $\text{var}[t] = E[t^2] - (E[t])^2 = 2/\lambda^2 - 1/\lambda^2 = 1/\lambda^2$. \square

• **h.** 2 points Which multiple of t has the lowest MSE as an estimator of $1/\lambda$?

Answer. It is $t/2$. Just plug $\theta = 1/\lambda$ into (24).

$$(26) \quad a = \frac{1/\lambda^2}{\text{var}[t] + 1/\lambda^2} = \frac{1/\lambda^2}{1/\lambda^2 + 1/\lambda^2} = \frac{1}{2}.$$

 \square

• **i.** 2 points Assume t_1, \dots, t_n are independently distributed, and each of them has the exponential distribution with the same parameter λ . Which multiple of the sample mean $\bar{t} = \frac{1}{n} \sum_{i=1}^n t_i$ has best MSE as estimator of $1/\lambda$?

Answer. \bar{t} has expected value $1/\lambda$ and variance $1/n\lambda^2$. Therefore

$$(27) \quad a = \frac{1/\lambda^2}{\text{var}[\bar{t}] + 1/\lambda^2} = \frac{1/\lambda^2}{1/n\lambda^2 + 1/\lambda^2} = \frac{n}{n+1},$$

i.e., for the best estimator $\tilde{t} = \frac{1}{n+1} \sum t_i$ divide the sum by $n+1$ instead of n . \square

• **j.** 3 points Assume $q \sim \sigma^2 \chi_m^2$ (in other words, $\frac{1}{\sigma^2} q \sim \chi_m^2$, 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 q that has minimum MSE as estimator of σ^2 .

Answer. This is a trick question since q itself is not an unbiased estimator of σ^2 . $E[q] = m\sigma^2$, therefore q/m is the unbiased estimator. Since $\text{var}[q/m] = 2\sigma^4/m$, it follows from (24) that $a = m/(m+2)$, therefore the minimum MSE multiple of q is $\frac{q}{m} \frac{m}{m+2} = \frac{q}{m+2}$. I.e., divide q by $m+2$ instead of m . \square

• **k.** 3 points Assume you have n independent observations of a Normally distributed random variable y with unknown mean μ and standard deviation σ^2 . The best unbiased estimator of σ^2 is $\frac{1}{n-1} \sum (y_i - \bar{y})^2$, and the maximum likelihood estimator is $\frac{1}{n} \sum (y_i - \bar{y})^2$. 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 $\sum (y_i - \bar{y})^2$?

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 $\sum (y_i - \bar{y})^2$ is $\frac{1}{n+1} \sum (y_i - \bar{y})^2$. \square

• **l.** 3 points We are still in the model defined in k. Which multiple of the sample mean \bar{y} has smallest MSE as estimator of μ ? 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{y} . Indeed, \bar{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 μ^2 and $\mu^2 + \frac{1}{n}\sigma^2$ from the data, shows that such an operationalization is sometimes successful. \square

Problem 27. 4 points Assume n independent observations of a variable $y \sim N(\mu, \sigma^2)$ are available, where σ^2 is known. Show that the sample mean \bar{y} attains the Cramer-Rao lower bound for μ .

Answer. The density function of each y_i is

$$(28) \quad f_{y_i}(y) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

therefore the log likelihood function of the whole vector is

$$(29) \quad \ell(\mathbf{y}; \mu) = \sum_{i=1}^n \log f_{y_i}(y_i) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2$$

$$(30) \quad \frac{\partial}{\partial \mu} \ell(\mathbf{y}; \mu) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu)$$

In order to apply (??) you can either square this and take the expected value

$$(31) \quad \mathbb{E}\left[\left(\frac{\partial}{\partial \mu} \ell(\mathbf{y}; \mu)\right)^2\right] = \frac{1}{\sigma^4} \sum \mathbb{E}[(y_i - \mu)^2] = n/\sigma^2$$

alternatively one may take one more derivative from (30) to get

$$(32) \quad \frac{\partial^2}{\partial \mu^2} \ell(\mathbf{y}; \mu) = -\frac{n}{\sigma^2}$$

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 28. *3 points* You have a Bernoulli experiment with unknown parameter θ , $0 \leq \theta \leq 1$. 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 (??) $\binom{12}{9} \theta^9 (1-\theta)^3$, and in the second it is by (??) $\binom{11}{8} \theta^9 (1-\theta)^3$. They are proportional, the stopping rule therefore does not matter! \square

Problem 29. *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 in 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. \square

- 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, i.e., it is the opposite way. \square

Problem 30. 1 point Compute the matrix product

$$\begin{bmatrix} 1 & 2 & 4 \\ 0 & 3 & 3 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 2 & 1 \\ 3 & 8 \end{bmatrix}$$

Answer.

$$\begin{bmatrix} 1 & 2 & 4 \\ 0 & 3 & 3 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 2 & 1 \\ 3 & 8 \end{bmatrix} = \begin{bmatrix} 1 \cdot 4 + 2 \cdot 2 + 4 \cdot 3 & 1 \cdot 0 + 2 \cdot 1 + 4 \cdot 8 \\ 0 \cdot 4 + 3 \cdot 2 + 3 \cdot 3 & 0 \cdot 0 + 3 \cdot 1 + 3 \cdot 8 \end{bmatrix} = \begin{bmatrix} 20 & 34 \\ 15 & 27 \end{bmatrix}$$

\square

Problem 31. 4 points Using the matrix differentiation rules

$$(33) \quad \partial \mathbf{w}^\top \mathbf{x} / \partial \mathbf{x}^\top = \mathbf{w}^\top$$

$$(34) \quad \partial \mathbf{x}^\top \mathbf{M} \mathbf{x} / \partial \mathbf{x}^\top = 2 \mathbf{x}^\top \mathbf{M}$$

for symmetric \mathbf{M} , compute the least-squares estimate $\hat{\beta}$ which minimizes

$$(35) \quad SSE = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta)$$

You are allowed to assume that $\mathbf{X}^\top \mathbf{X}$ has an inverse.

Answer. First you have to multiply out

$$(36) \quad (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}.$$

The matrix differentiation rules (33) and (34) allow us to differentiate (36) to get

$$(37) \quad \partial SSE / \partial \boldsymbol{\beta}^\top = -2\mathbf{y}^\top \mathbf{X} + 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}.$$

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”

$$(38) \quad \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}}.$$

Since $\mathbf{X}^\top \mathbf{X}$ has an inverse, one can premultiply both sides of (38) by $(\mathbf{X}^\top \mathbf{X})^{-1}$:

$$(39) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

□

Problem 32. 1 point In the decomposition

$$(40) \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\epsilon}}$$

which of $\boldsymbol{\epsilon}$ and $\hat{\boldsymbol{\epsilon}}$ is called the residual and which is called the disturbance? Which of \mathbf{y} , \mathbf{X} , $\boldsymbol{\beta}$, $\boldsymbol{\epsilon}$, $\hat{\boldsymbol{\beta}}$, and $\hat{\boldsymbol{\epsilon}}$ is observed (or a function of observed quantities) and which is unobserved?

Answer. $\boldsymbol{\epsilon}$ is called the disturbance and $\hat{\boldsymbol{\epsilon}}$ the residual. $\boldsymbol{\beta}$ and $\boldsymbol{\epsilon}$ are unobserved, all others are observed.

□

Problem 33. 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, $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{y}$. Idempotent, i.e. $\mathbf{M}\mathbf{M} = \mathbf{M}$:

(41)

$$\mathbf{M}\mathbf{M} = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{M}$$

□

Problem 34. 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 $\mathbf{X}^\top \mathbf{q} = \mathbf{o}$. □

• **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 $\mathbf{M}\mathbf{q} = \mathbf{M}\mathbf{X}\mathbf{a} = \mathbf{O}\mathbf{a} = \mathbf{o}$. □

Problem 35. 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 $\mathcal{C}[\hat{\mathbf{y}}, \hat{\boldsymbol{\varepsilon}}] = \mathbf{O}$. Here the covariance matrix $\mathcal{C}[\mathbf{x}, \mathbf{z}]$ is that matrix whose (i, j) element is $\text{cov}[x_i, z_j]$. The covariance matrix satisfies the rules $\mathcal{C}[\mathbf{B}\mathbf{y}, \mathbf{T}\mathbf{z}] = \mathbf{B}\mathcal{C}[\mathbf{y}, \mathbf{z}]\mathbf{T}^\top$, $\mathcal{C}[\mathbf{y}, \mathbf{y}] = \mathcal{V}[\mathbf{y}]$, $\mathcal{C}[\mathbf{z}, \mathbf{y}] = (\mathcal{C}[\mathbf{y}, \mathbf{z}])^\top$, $\mathcal{C}[\mathbf{x} + \mathbf{y}, \mathbf{z}] = \mathcal{C}[\mathbf{x}, \mathbf{z}] + \mathcal{C}[\mathbf{y}, \mathbf{z}]$, and $\mathcal{C}[\mathbf{x}, \mathbf{c}] = \mathbf{O}$ if \mathbf{c} is a vector of constants.

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 ??.

□

-
Maximum number of points: 75.

REFERENCES

- [KS79] Sir Maurice Kendall and Alan Stuart, *The advanced theory of statistics*, fourth ed., vol. 2, Griffin, London, 1979. 7

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

E-mail address: `ehrbar@econ.utah.edu`

URL: `http://www.econ.utah.edu/ehrbar`