

Chapter 3: Multiple Regression Analysis: Estimation Problems
 Book Title: Introductory Econometrics
 Printed By: Wanwiphang Manachotipong (wanwiphang@econ.tu.ac.th)
 © 2016 Cengage Learning, Cengage Learning

Chapter Review

Problems

1. Using the data in GPA2 on 4,137 college students, the following equation was estimated by OLS:

$$\widehat{colgpa} = 1.392 - .0135 \text{ hsperc} + .00148 \text{ sat}$$

$$n = 4,137, R^2 = .273,$$

where *colgpa* is measured on a four-point scale, *hsperc* is the percentile in the high school graduating class (defined so that, for example, *hsperc* = 5 means the top 5% of the class), and *sat* is the combined math and verbal scores on the student achievement test.

- i. Why does it make sense for the coefficient on *hsperc* to be negative?
 - ii. What is the predicted college GPA when *hsperc* = 20 and *sat* = 1,050?
 - iii. Suppose that two high school graduates, A and B, graduated in the same percentile from high school, but Student A's SAT score was 140 points higher (about one standard deviation in the sample). What is the predicted difference in college GPA for these two students? Is the difference large?
 - iv. Holding *hsperc* fixed, what difference in SAT scores leads to a predicted *colgpa* difference of .50, or one-half of a grade point? Comment on your answer.
2. The data in WAGE2 on working men was used to estimate the following equation:

$$\widehat{educ} = 10.36 - .094 \text{ sibs} + .131 \text{ meduc} + .210 \text{ feduc}$$

$$n = 722, R^2 = .214,$$

where *educ* is years of schooling, *sibs* is number of siblings, *meduc* is mother's years of schooling, and *feduc* is father's years of schooling.

- i. Does *sibs* have the expected effect? Explain. Holding *meduc* and *feduc* fixed, by how much does *sibs* have to increase to reduce predicted years of education by one year? (A noninteger answer is acceptable here.)
 - ii. Discuss the interpretation of the coefficient on *meduc*.
 - iii. Suppose that Man A has no siblings, and his mother and father each have 12 years of education. Man B has no siblings, and his mother and father each have 16 years of education. What is the predicted difference in years of education between B and A?
3. The following model is a simplified version of the multiple regression model used by Biddle and Hamermesh (1990) to study the tradeoff between time spent sleeping and working and to look at other factors affecting sleep:

$$\text{sleep} = \beta_0 + \beta_1 \text{totwrk} + \beta_2 \text{educ} + \beta_3 \text{age} + u,$$

where *sleep* and *totwrk* (total work) are measured in minutes per week and *educ* and *age* are measured in years. (See also [Computer Exercise C3](#) in [Chapter 2](#).)

- i. If adults trade off sleep for work, what is the sign of β_1 ?
- ii. What signs do you think β_2 and β_3 will have?
- iii. Using the data in SLEEP75, the estimated equation is

$$\widehat{\text{sleep}} = 3,638.25 - .148 \text{totwrk} - 11.13 \text{educ} + 2.20 \text{age}$$

$$n = 706, R^2 = .113.$$

If someone works five more hours per week, by how many minutes is *sleep* predicted to fall? Is this a large tradeoff?

- iv. Discuss the sign and magnitude of the estimated coefficient on *educ*.
 - v. Would you say *totwrk*, *educ*, and *age* explain much of the variation in *sleep*? What other factors might affect the time spent sleeping? Are these likely to be correlated with *totwrk*?
4. The median starting salary for new law school graduates is determined by

$$\log(\text{salary}) = \beta_0 + \beta_1 \text{LSAT} + \beta_2 \text{GPA} + \beta_3 \log(\text{libvol}) + \beta_4 \log(\text{cost}) + \beta_5 \text{rank} + u,$$

where *LSAT* is the median LSAT score for the graduating class, *GPA* is the median college GPA for the class, *libvol* is the number of volumes in the law

school library, $cost$ is the annual cost of attending law school, and $rank$ is a law school ranking (with $rank = 1$ being the best).

- i. Explain why we expect $\beta_5 \leq 0$.
- ii. What signs do you expect for the other slope parameters? Justify your answers.
- iii. Using the data in LAWSCH85, the estimated equation is

$$\widehat{\log(salary)} = 8.34 + .0047 LAST + .248 GPA + .095 \log(libvol) + .038 \log(cost) - .0033 rank$$

$$n = 136, R^2 = .842.$$

What is the predicted ceteris paribus difference in salary for schools with a median GPA different by one point? (Report your answer as a percentage.)

- iv. Interpret the coefficient on the variable $\log(libvol)$.
 - v. Would you say it is better to attend a higher ranked law school? How much is a difference in ranking of 20 worth in terms of predicted starting salary?
5. In a study relating college grade point average to time spent in various activities, you distribute a survey to several students. The students are asked how many hours they spend each week in four activities: studying, sleeping, working, and leisure. Any activity is put into one of the four categories, so that for each student, the sum of hours in the four activities must be 168.

- i. In the model

$$GPA = \beta_0 + \beta_1 study + \beta_2 sleep + \beta_3 work + \beta_4 leisure + u,$$

does it make sense to hold $sleep$, $work$, and $leisure$ fixed, while changing $study$?

- ii. Explain why this model violates [Assumption MLR.3](#).
 - iii. How could you reformulate the model so that its parameters have a useful interpretation and it satisfies [Assumption MLR.3](#)?
6. Consider the multiple regression model containing three independent variables, under [Assumptions MLR.1](#), [MLR.2](#), [MLR.3](#) and [MLR.4](#):

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u.$$

You are interested in estimating the sum of the parameters on x_1 and x_2 ; call this $\theta_1 = \beta_0 + \beta_1$.

- i. Show that $\hat{\theta}_1 = \hat{\beta}_1 + \hat{\beta}_2$ is an unbiased estimator of θ_1 .
- ii. Find $\text{Var}(\hat{\theta}_1)$ in terms of $\text{Var}(\hat{\beta}_1)$, $\text{Var}(\hat{\beta}_2)$, and $\text{Corr}(\hat{\beta}_1, \hat{\beta}_2)$.

7. Which of the following can cause OLS estimators to be biased?

- i. Heteroskedasticity.
- ii. Omitting an important variable.
- iii. A sample correlation coefficient of .95 between two independent variables both included in the model.

8. Suppose that average worker productivity at manufacturing firms (*avgprod*) depends on two factors, average hours of training (*avgtrain*) and average worker ability (*avgabil*):

$$\text{avgprod} = \beta_0 + \beta_1 \text{avgtrain} + \beta_2 \text{avgabil} + u.$$

Assume that this equation satisfies the Gauss-Markov assumptions. If grants have been given to firms whose workers have less than average ability, so that *avgtrain* and *avgabil* are negatively correlated, what is the likely bias in $\tilde{\beta}_1$ obtained from the simple regression of *avgprod* on *avgtrain*?

9. The following equation describes the median housing price in a community in terms of amount of pollution (*nox* for nitrous oxide) and the average number of rooms in houses in the community (*rooms*):

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{nox}) + \beta_2 \text{rooms} + u.$$

- i. What are the probable signs of β_1 and β_2 ? What is the interpretation of β_1 ? Explain.
- ii. Why might *nox* [or more precisely, $\log(\text{nox})$] and *rooms* be negatively correlated? If this is the case, does the simple regression of $\log(\text{price})$ on $\log(\text{nox})$ produce an upward or a downward biased estimator of β_1 ?
- iii. Using the data in HPRICE2, the following equations were estimated:

$$\widehat{\log(\text{price})} = 11.71 - 1.043 \log(\text{nox}), n = 506, R^2 = .264.$$

$$\widehat{\log(\text{price})} = 9.23 - .718 \log(\text{nox}) + .306 \text{rooms}, n = 506, R^2 = .514.$$

Is the relationship between the simple and multiple regression estimates of the elasticity of *price* with respect to *nox* what you would have predicted, given your answer in part (ii)? Does this mean that $-.718$ is definitely closer to the true elasticity than -1.043 ?

10. Suppose that you are interested in estimating the ceteris paribus relationship between y and x_1 . For this purpose, you can collect data on two control variables, x_2 and x_3 . (For concreteness, you might think of y as final exam score, x_1 as class attendance, x_2 as GPA up through the previous semester, and x_3 as SAT or ACT score.) Let $\tilde{\beta}_1$ be the simple regression estimate from y on x_1 and let $\hat{\beta}_1$ be the multiple regression estimate from y on x_1, x_2, x_3 .
- If x_1 is highly correlated with x_2 and x_3 in the sample, and x_2 and x_3 have large partial effects on y , would you expect $\tilde{\beta}_1$ and $\hat{\beta}_1$ to be similar or very different? Explain.
 - If x_1 is almost uncorrelated with x_2 and x_3 , but x_2 and x_3 are highly correlated, will $\tilde{\beta}_1$ and $\hat{\beta}_1$ tend to be similar or very different? Explain.
 - If x_1 is highly correlated with x_2 and x_3 , and x_2 and x_3 have small partial effects on y , would you expect $\text{se}(\tilde{\beta}_1)$ or $\text{se}(\hat{\beta}_1)$ to be smaller? Explain.
 - If x_1 is almost uncorrelated with x_2 and x_3 , x_2 and x_3 have large partial effects on y , and x_2 and x_3 are highly correlated, would you expect $\text{se}(\tilde{\beta}_1)$ or $\text{se}(\hat{\beta}_1)$ to be smaller? Explain.

11. Suppose that the population model determining y is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u,$$

and this model satisfies [Assumptions MLR.1](#), [MLR.2](#), [MLR.3](#) and [MLR.4](#). However, we estimate the model that omits x_3 . Let $\tilde{\beta}_0$, $\tilde{\beta}_1$, and $\tilde{\beta}_2$ be the OLS estimators from the regression of y on x_1 and x_2 . Show that the expected value of $\tilde{\beta}_1$ (given the values of the independent variables in the sample) is

$$\mathbb{E}(\tilde{\beta}_1) = \beta_1 + \beta_3 \frac{\sum_{i=1}^n \hat{r}_{i1} x_{i3}}{\sum_{i=1}^n \hat{r}_{i1}^2},$$

where the \hat{r}_{i1} are the OLS residuals from the regression of x_1 on x_2 . [*Hint:* The formula for $\tilde{\beta}_1$ comes from [equation \(3.22\)](#). Plug $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + u_i$ into this equation. After some algebra, take the expectation treating x_{i3} and \hat{r}_{i1} as nonrandom.]

12. The following equation represents the effects of tax revenue mix on subsequent employment growth for the population of counties in the United States:

$$growth = \beta_0 + \beta_1 share_P + \beta_2 share_I + \beta_3 share_S + other\ factors,$$

where *growth* is the percentage change in employment from 1980 to 1990, *share_P* is the share of property taxes in total tax revenue, *share_I* is the share of income tax revenues, and *share_S* is the share of sales tax revenues. All of these variables are measured in 1980. The omitted share, *share_F*, includes fees and miscellaneous taxes. By definition, the four shares add up to one. Other factors would include expenditures on education, infrastructure, and so on (all measured in 1980).

- i. Why must we omit one of the tax share variables from the equation?
 - ii. Give a careful interpretation of β_1 .
13. i. Consider the simple regression model $y = \beta_0 + \beta_1 x + u$ under the first four Gauss-Markov assumptions. For some function $g(x)$, for example $g(x) = x^2$ or $g(x) = \log(1 + x^2)$, define $z_i = g(x_i)$. Define a slope estimator as

$$\tilde{\beta}_1 = \left(\sum_{i=1}^n (z_i - \bar{z}) y_i \right) / \left(\sum_{i=1}^n (z_i - \bar{z}) x_i \right).$$

Show that $\tilde{\beta}_1$ is linear and unbiased. Remember, because $\mathbf{E}(u|x) = 0$, you can treat both x_i and z_i as nonrandom in your derivation.

- ii. Add the homoskedasticity assumption, [MLR.5](#). Show that

$$\text{Var}(\tilde{\beta}_1) = \sigma^2 \left(\sum_{i=1}^n (z_i - \bar{z})^2 \right) / \left(\sum_{i=1}^n (z_i - \bar{z}) x_i \right)^2.$$

- iii. Show directly that, under the Gauss-Markov assumptions, $\text{Var}(\hat{\beta}_1) \leq \text{Var}(\tilde{\beta}_1)$, where $\hat{\beta}_1$ is the OLS estimator. [*Hint*: The Cauchy-Schwartz inequality in [Appendix B](#) implies that

$$\left(n^{-1} \sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x}) \right)^2 \leq \left(n^{-1} \sum_{i=1}^n (z_i - \bar{z})^2 \right) \left(n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right);$$

notice that we can drop \bar{x} from the sample covariance.]

14. Suppose you have a sample of size n on three variables, y , x_1 , and x_2 , and you are primarily interested in the effect of x_1 on y . Let $\tilde{\beta}_1$ be the coefficient on x_1 from the simple regression and $\hat{\beta}_1$ the coefficient on x_1 from the regression y on x_1, x_2 . The standard errors reported by any regression package are

$$se(\tilde{\beta}_1) = \frac{\tilde{\sigma}}{\sqrt{SST_1}}$$

$$se(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{SST_1}} \cdot \sqrt{VIF_1},$$

where $\tilde{\sigma}$ is the SER from the simple regression, $\hat{\sigma}$ is the SER from the multiple regression, $VIF_1 = 1/(1 - R_1^2)$, and R_1^2 is the R -squared from the regression of x_1 on x_2 . Explain why $se(\hat{\beta}_1)$ can be smaller or larger than $se(\tilde{\beta}_1)$.

15. The following estimated equations use the data in MLB1, which contains information on major league baseball salaries. The dependent variable, $\ln salary$, is the log of salary. The two explanatory variables are years in the major leagues ($years$) and runs batted in per year ($rbisyr$):

$$\widehat{\ln salary} = 12.373 + .1770 \text{ years}$$

$$(.098) \quad (.0132)$$

$$n = 353, SSR = 326.196, SER = .964, R^2 = .337$$

$$\widehat{\ln salary} = 11.861 + .0904 \text{ years} + .0302 \text{ rbisyr}$$

$$(.084) \quad (.0118) \quad (.0020)$$

$$n = 353, SSR = 198.475, SER = .753, R^2 = .597$$

- i. How many degrees of freedom are in each regression? How come the SER is smaller in the second regression than the first?
 - ii. The sample correlation coefficient between $years$ and $rbisyr$ is about 0.487. Does this make sense? What is the variance inflation factor (there is only one) for the slope coefficients in the multiple regression? Would you say there is little, moderate, or strong collinearity between $years$ and $rbisyr$?
 - iii. How come the standard error for the coefficient on $years$ in the multiple regression is lower than its counterpart in the simple regression?
16. The following equations were estimated using the data in LAWSCH85:

$$\widehat{\ln salary} = 9.90 - .0041 \text{ rank} + .294 \text{ GPA}$$

$$(.24) \quad (.0003) \quad (.069)$$

$$n = 142, R^2 = .8238$$

$$\widehat{\ln salary} = 9.86 - .0038 \text{ rank} + .295 \text{ GPA} + .00017 \text{ age}$$

$$(.29) \quad (.0004) \quad (.083) \quad (.00036)$$

$$n = 99, R^2 = .8036$$

How can it be that the R-squared is smaller when the variable *age* is added to the equation?

Chapter 3: Multiple Regression Analysis: Estimation Problems

Book Title: Introductory Econometrics

Printed By: Wanwiphang Manachotipong (wanwiphang@econ.tu.ac.th)

© 2016 Cengage Learning, Cengage Learning

© 2020 Cengage Learning Inc. All rights reserved. No part of this work may be reproduced or used in any form or by any means - graphic, electronic, or mechanical, or in any other manner - without the written permission of the copyright holder.