

Chapter Review

Problems

1. The following equation was estimated using the data in CEOSAL1:

$$\widehat{\log(\text{salary})} = 4.322 + .276 \log(\text{sales}) + .0215 \text{roe} - .00008 \text{roe}^2$$

$$\begin{array}{cccc} (.324) & (.033) & (.0129) & (.00026) \end{array}$$

$$n = 209, R^2 = .282.$$

This equation allows *roe* to have a diminishing effect on $\log(\text{salary})$. Is this generality necessary? Explain why or why not.

2. Let $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ be the OLS estimates from the regression of y_i on x_{i1}, \dots, x_{ik} , $i = 1, 2, \dots, n$. For nonzero constants c_1, \dots, c_k , argue that the OLS intercept and slopes from the regression of $c_0 y_i$ on $c_1 x_{i1}, \dots, c_k x_{ik}$, $i = 1, 2, \dots, n$, are given by $\tilde{\beta}_0 = c_0 \hat{\beta}_0$, $\tilde{\beta}_1 = (c_0/c_1) \hat{\beta}_1, \dots, \tilde{\beta}_k = (c_0/c_k) \hat{\beta}_k$. [Hint: Use the fact that the $\hat{\beta}_j$ solve the first order conditions in (3.13), and the $\tilde{\beta}_j$ must solve the first order conditions involving the rescaled dependent and independent variables.]
3. Using the data in RDCHEM, the following equation was obtained by OLS:

$$\widehat{\text{rdintens}} = 2.613 + .00030 \text{sales} - .0000000070 \text{sales}^2$$

$$\begin{array}{ccc} (.429) & (.00014) & (.0000000037) \end{array}$$

$$n = 32, R^2 = .1484.$$

- i. At what point does the marginal effect of *sales* on *rdintens* become negative?
- ii. Would you keep the quadratic term in the model? Explain.
- iii. Define *salesbil* as sales measured in billions of dollars:
 $\text{salesbil} = \text{sales}/1,000$. Rewrite the estimated equation with *salesbil* and salesbil^2 as the independent variables. Be sure to report standard errors and the *R*-squared. [Hint: Note that $\text{salesbil}^2 = \text{sales}^2/(1,000)^2$.]
- iv. For the purpose of reporting the results, which equation do you prefer?

4. The following model allows the return to education to depend upon the total amount of both parents' education, called *pareduc*:

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{educ} \cdot \text{pareduc} + \beta_3 \text{exper} + \beta_4 \text{tenure} + u.$$

- i. Show that, in decimal form, the return to another year of education in this model is

$$\Delta \log(\text{wage}) / \Delta \text{educ} = \beta_1 + \beta_2 \text{pareduc}.$$

What sign do you expect for β_2 ? Why?

- ii. Using the data in WAGE2, the estimated equation is

$$\begin{aligned} \widehat{\log(\text{wage})} &= 5.65 + .047 \text{educ} + .00078 \text{educ} \cdot \text{pareduc} + \\ &\quad (.13) \quad (.010) \quad (.00021) \\ &\quad .019 \text{exper} + .010 \text{tenure} \\ &\quad (.004) \quad (.003) \\ n &= 722, R^2 = .169. \end{aligned}$$

(Only 722 observations contain full information on parents' education.) Interpret the coefficient on the interaction term. It might help to choose two specific values for *pareduc*—for example, *pareduc* = 32 if both parents have a college education, or *pareduc* = 24 if both parents have a high school education—and to compare the estimated return to *educ*.

- iii. When *pareduc* is added as a separate variable to the equation, we get:

$$\begin{aligned} \widehat{\log(\text{wage})} &= 4.94 + .097 \text{educ} + .033 \text{pareduc} - .0016 \text{educ} \cdot \text{pareduc} \\ &\quad (.38) \quad (.027) \quad (.017) \quad (.0012) \\ &\quad + .020 \text{exper} + .010 \text{tenure} \\ &\quad (.004) \quad (.003) \\ n &= 722, R^2 = .174. \end{aligned}$$

Does the estimated return to education now depend positively on parent education? Test the null hypothesis that the return to education does not depend on parent education.

5. In [Example 4.2](#), where the percentage of students receiving a passing score on a tenth-grade math exam (*math10*) is the dependent variable, does it make sense to include *sci11*—the percentage of eleventh graders passing a science exam—as an additional explanatory variable?
6. When *atndrte*² and *ACT*·*atndrte* are added to the equation estimated in [\(6.19\)](#), the *R*-squared becomes .232. Are these additional terms jointly

significant at the 10% level? Would you include them in the model?

7. The following three equations were estimated using the 1,534 observations in 401K:

$$\widehat{prate} = 80.29 + 5.44 \text{ mrate} + .269 \text{ age} - .00013 \text{ totemp}$$

(.78)
(.52)
(.045)
(.00004)

$$R^2 = .100, \bar{R}^2 = .098.$$

$$\widehat{prate} = 97.32 + 5.02 \text{ mrate} + .314 \text{ age} - 2.66 \log(\text{totemp})$$

(1.95)
(0.51)
(.044)
(.28)

$$R^2 = .144, \bar{R}^2 = .142.$$

$$\widehat{prate} = 80.62 + 5.34 \text{ mrate} + .290 \text{ age} - .00043 \text{ totemp} + .0000000039 \text{ totemp}^2$$

(.78)
(.52)
(.045)
(.00009)
(.00000000010)

$$R^2 = .108, \bar{R}^2 = .106.$$

Which of these three models do you prefer? Why?

8. Suppose we want to estimate the effects of alcohol consumption (*alcohol*) on college grade point average (*colGPA*). In addition to collecting information on grade point averages and alcohol usage, we also obtain attendance information (say, percentage of lectures attended, called *attend*). A standardized test score (say, *SAT*) and high school GPA (*hsGPA*) are also available.

i. Should we include *attend* along with *alcohol* as explanatory variables in a multiple regression model? (Think about how you would interpret β_{alcohol} .)

ii. Should *SAT* and *hsGPA* be included as explanatory variables? Explain.

9. If we start with (6.38) under the CLM assumptions, assume large n , and ignore the estimation error in the $\hat{\beta}_j$, a 95% prediction interval for y^0 is $[\exp(-1.96\hat{\sigma}) \exp(\widehat{\log y^0}), \exp(1.96\hat{\sigma}) \exp(\widehat{\log y^0})]$. The point prediction for y^0 is $\hat{y}^0 = \exp(\hat{\sigma}_2) \exp(\widehat{\log y^0})$.

i. For what values of $\hat{\sigma}$ will the point prediction be in the 95% prediction interval? Does this condition seem likely to hold in most applications?

ii. Verify that the condition from part (i) is satisfied in the CEO salary example.

10. The following two equations were estimated using the data in MEAPSINGLE. The key explanatory variable is *lexppp*, the log of expenditures per student at the school level.

$$\widehat{math4} = 24.49 + 9.01 lexppp - .422 free - .752 lmedinc - .274 pctsgle$$

$$(59.24) (4.04) (.071) (5.358) (.161)$$

$$n = 229, R^2 = .472, \bar{R}^2 = .462.$$

$$\widehat{math4} = 149.38 + 1.93 lexppp - .060 free - 10.78 lmedinc - .397 pctsgle + .667 read4$$

$$(41.70) (2.82) (.054) (3.76) (.111) (.042)$$

$$n = 229, R^2 = .749, \bar{R}^2 = .743.$$

- i. If you are a policy maker trying to estimate the causal effect of per-student spending on math test performance, explain why the first equation is more relevant than the second. What is the estimated effect of a 10% increase in expenditures per student?
- ii. Does adding *read4* to the regression have strange effects on coefficients and statistical significance other than β_{lexppp} ?
- iii. How would you explain to someone with only basic knowledge of regression why, in this case, you prefer the equation with the smaller adjusted *R*-squared?

Chapter 6: Multiple Regression Analysis: Further Issues 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.