

PROBLEMS

- 6.1** 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 *scill*—the percentage of eleventh graders passing a science exam—as an additional explanatory variable?
- 6.2** 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 $\hat{y}^0 = \exp(\hat{\sigma}^2/2) \exp(\widehat{\log y^0})$.
- (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.
- 6.3** 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.RAW, 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.

4.4 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.

4.5 Using the data in RDCHEM.RAW, the following equation was obtained by OLS:

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

$$(.429) \quad (.00014) \quad (.0000000037)$$

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

- (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?

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

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

$$(.78) \quad (.52) \quad (.045) \quad (.00004)$$

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

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

$$(1.95) \quad (0.51) \quad (.044) \quad (.28)$$

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

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

$$(.78) \quad (.52) \quad (.045) \quad (.00009)$$

$$+ .0000000039 \text{ totemp}^2$$

$$(.0000000010)$$

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

Which of these three models do you prefer? Why?

4.7 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.]

- 6.8** When $atndrte^2$ and $ACT \cdot atndrte$ 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?
- 6.9** The following equation was estimated using the data in CEOSAL1.RAW:

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

$$(.324) \quad (.033) \quad (.0129) \quad (.00026)$$

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

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

COMPUTER EXERCISES

- C6.1** Use the data in KIELMC.RAW, only for the year 1981, to answer the following questions. The data are for houses that sold during 1981 in North Andover, Massachusetts; 1981 was the year construction began on a local garbage incinerator.

- (i) To study the effects of the incinerator location on housing price, consider the simple regression model

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{dist}) + u,$$

where price is housing price in dollars and dist is distance from the house to the incinerator measured in feet. Interpreting this equation causally, what sign do you expect for β_1 if the presence of the incinerator depresses housing prices? Estimate this equation and interpret the results.

- (ii) To the simple regression model in part (i), add the variables $\log(\text{intst})$, $\log(\text{area})$, $\log(\text{land})$, rooms , baths , and age , where intst is distance from the home to the interstate, area is square footage of the house, land is the lot size in square feet, rooms is total number of rooms, baths is number of bathrooms, and age is age of the house in years. Now, what do you conclude about the effects of the incinerator? Explain why (i) and (ii) give conflicting results.
- (iii) Add $[\log(\text{intst})]^2$ to the model from part (ii). Now what happens? What do you conclude about the importance of functional form?
- (iv) Is the square of $\log(\text{dist})$ significant when you add it to the model from part (iii)?
- C6.2** Use the data in WAGE1.RAW for this exercise.
- (i) Use OLS to estimate the equation

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{exper}^2 + u$$

and report the results using the usual format.

- (ii) Is exper^2 statistically significant at the 1% level?

- (iii) Using the approximation

$$\% \Delta \widehat{wage} \approx 100(\hat{\beta}_2 + 2\hat{\beta}_3 \text{exper}) \Delta \text{exper},$$

find the approximate return to the fifth year of experience. What is the approximate return to the twentieth year of experience?

- (iv) At what value of *exper* does additional experience actually lower predicted $\log(\text{wage})$? How many people have more experience in this sample?

C6.3 Consider a model where the return to education depends upon the amount of work experience (and vice versa):

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

- (i) Show that the return to another year of education (in decimal form), holding *exper* fixed, is $\beta_1 + \beta_3 \text{exper}$.
- (ii) State the null hypothesis that the return to education does not depend on the level of *exper*. What do you think is the appropriate alternative?
- (iii) Use the data in WAGE2.RAW to test the null hypothesis in (ii) against your stated alternative.
- (iv) Let θ_1 denote the return to education (in decimal form), when *exper* = 10: $\theta_1 = \beta_1 + 10\beta_3$. Obtain $\hat{\theta}_1$ and a 95% confidence interval for θ_1 . (Hint: Write $\beta_1 = \theta_1 - 10\beta_3$ and plug this into the equation; then rearrange. This gives the regression for obtaining the confidence interval for θ_1 .)

C6.4 Use the data in GPA2.RAW for this exercise.

- (i) Estimate the model

$$\text{sat} = \beta_0 + \beta_1 \text{hsize} + \beta_2 \text{hsize}^2 + u,$$

where *hsize* is the size of the graduating class (in hundreds), and write the results in the usual form. Is the quadratic term statistically significant?

- (ii) Using the estimated equation from part (i), what is the “optimal” high school size? Justify your answer.
- (iii) Is this analysis representative of the academic performance of *all* high school seniors? Explain.
- (iv) Find the estimated optimal high school size, using $\log(\text{sat})$ as the dependent variable. Is it much different from what you obtained in part (ii)?

C6.5 Use the housing price data in HPRICE1.RAW for this exercise.

- (i) Estimate the model

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{lotsize}) + \beta_2 \log(\text{sqrft}) + \beta_3 \text{bdrms} + u$$

and report the results in the usual OLS format.

- (ii) Find the predicted value of $\log(\text{price})$, when $\text{lotsize} = 20,000$, $\text{sqrft} = 2,500$, and $\text{bdrms} = 4$. Using the methods in Section 6.4, find the predicted value of price at the same values of the explanatory variables.
- (iii) For explaining variation in price , decide whether you prefer the model from part (i) or the model

$$\text{price} = \beta_0 + \beta_1 \text{lotsize} + \beta_2 \text{sqrft} + \beta_3 \text{bdrms} + u.$$

C6.6 Use the data in VOTE1.RAW for this exercise.

- (i) Consider a model with an interaction between expenditures:

$$\text{voteA} = \beta_0 + \beta_1 \text{prtystrA} + \beta_2 \text{expendA} + \beta_3 \text{expendB} + \beta_4 \text{expendA} \cdot \text{expendB} + u.$$

What is the partial effect of expendB on voteA , holding prtystrA and expendA fixed? What is the partial effect of expendA on voteA ? Is the expected sign for β_4 obvious?

- (ii) Estimate the equation in part (i) and report the results in the usual form. Is the interaction term statistically significant?
- (iii) Find the average of expendA in the sample. Fix expendA at 300 (for \$300,000). What is the estimated effect of another \$100,000 spent by Candidate B on voteA ? Is this a large effect?
- (iv) Now fix expendB at 100. What is the estimated effect of $\Delta \text{expendA} = 100$ on voteA ? Does this make sense?
- (v) Now, estimate a model that replaces the interaction with shareA , Candidate A's percentage share of total campaign expenditures. Does it make sense to hold both expendA and expendB fixed, while changing shareA ?
- (vi) (Requires calculus) In the model from part (v), find the partial effect of expendB on voteA , holding prtystrA and expendA fixed. Evaluate this at $\text{expendA} = 300$ and $\text{expendB} = 0$ and comment on the results.

C6.7 Use the data in ATTEND.RAW for this exercise.

- (i) In the model of Example 6.3, argue that

$$\Delta \text{stndfnl} / \Delta \text{priGPA} \approx \beta_2 + 2\beta_4 \text{priGPA} + \beta_6 \text{atndrte}.$$

Use equation (6.19) to estimate the partial effect when $\text{priGPA} = 2.59$ and $\text{atndrte} = 82$. Interpret your estimate.

- (ii) Show that the equation can be written as

$$\begin{aligned} \text{stndfnl} = & \theta_0 + \beta_1 \text{atndrte} + \theta_2 \text{priGPA} + \beta_3 \text{ACT} + \beta_4 (\text{priGPA} - 2.59)^2 \\ & + \beta_5 \text{ACT}^2 + \beta_6 \text{priGPA} (\text{atndrte} - 82) + u, \end{aligned}$$

where $\theta_2 = \beta_2 + 2\beta_4(2.59) + \beta_6(82)$. (Note that the intercept has changed, but this is unimportant.) Use this to obtain the standard error of $\hat{\theta}_2$ from part (i).

- (iii) Suppose that, in place of $\text{priGPA}(\text{atndrte} - 82)$, you put $(\text{priGPA} - 2.59) \cdot (\text{atndrte} - 82)$. Now how do you interpret the coefficients on atndrte and priGPA ?

C6.8 Use the data in HPRICE1.RAW for this exercise.

- (i) Estimate the model

$$price = \beta_0 + \beta_1 lotsize + \beta_2 sqrft + \beta_3 bdrms + u$$

and report the results in the usual form, including the standard error of the regression. Obtain predicted price, when we plug in $lotsize = 10,000$, $sqrft = 2,300$, and $bdrms = 4$; round this price to the nearest dollar.

- (ii) Run a regression that allows you to put a 95% confidence interval around the predicted value in part (i). Note that your prediction will differ somewhat due to rounding error.
- (iii) Let $price^0$ be the unknown future selling price of the house with the characteristics used in parts (i) and (ii). Find a 95% CI for $price^0$ and comment on the width of this confidence interval.

C6.9 The data set NBASAL.RAW contains salary information and career statistics for 269 players in the National Basketball Association (NBA).

- (i) Estimate a model relating points-per-game ($points$) to years in the league ($exper$), age , and years played in college ($coll$). Include a quadratic in $exper$; the other variables should appear in level form. Report the results in the usual way.
- (ii) Holding college years and age fixed, at what value of experience does the next year of experience actually reduce points-per-game? Does this make sense?
- (iii) Why do you think $coll$ has a negative and statistically significant coefficient? (*Hint*: NBA players can be drafted before finishing their college careers and even directly out of high school.)
- (iv) Add a quadratic in age to the equation. Is it needed? What does this appear to imply about the effects of age, once experience and education are controlled for?
- (v) Now regress $\log(wage)$ on $points$, $exper$, $exper^2$, age , and $coll$. Report the results in the usual format.
- (vi) Test whether age and $coll$ are jointly significant in the regression from part (v). What does this imply about whether age and education have separate effects on wage, once productivity and seniority are accounted for?

C6.10 Use the data in BWGHT2.RAW for this exercise.

- (i) Estimate the equation

$$\log(bwght) = \beta_0 + \beta_1 npvis + \beta_2 npvis^2 + u$$

by OLS, and report the results in the usual way. Is the quadratic term significant?

- (ii) Show that, based on the equation from part (i), the number of prenatal visits that maximizes $\log(bwght)$ is estimated to be about 22. How many women had at least 22 prenatal visits in the sample?
- (iii) Does it make sense that birth weight is actually predicted to decline after 22 prenatal visits? Explain.
- (iv) Add mother's age to the equation, using a quadratic functional form. Holding $npvis$ fixed, at what mother's age is the birth weight of the child maximized? What fraction of women in the sample are older than the "optimal" age?