

Further Inference in the Multiple Regression Model

Chapter 6

Prepared by Vera Tabakova, East Carolina University

Chapter 6:

Further Inference in the Multiple Regression Model

- 6.1 The F -Test
- 6.2 Testing the Significance of the Model
- 6.3 An Extended Model
- 6.4 Testing Some Economic Hypotheses
- 6.5 The Use of Nonsample Information
- 6.6 Model Specification
- 6.7 Poor Data, Collinearity and Insignificance
- 6.8 Prediction

6.1 The *F*-Test

$$S_i = \beta_1 + \beta_2 P_i + \beta_3 A_i + e_i \quad (6.1)$$

$$H_0 : \beta_2 = 0$$

$$H_1 : \beta_2 \neq 0$$

$$S_i = \beta_1 + \beta_3 A_i + e_i \quad (6.2)$$

6.1 The F -Test

$$F = \frac{(SSE_R - SSE_U) / J}{SSE_U / (N - K)} \quad (6.3)$$

If the null hypothesis is not true, then the difference between SSE_R and SSE_U becomes large, implying that the constraints placed on the model by the null hypothesis have a large effect on the ability of the model to fit the data.

6.1 The *F*-Test

Hypothesis testing steps:

1. Specify the null and alternative hypotheses: $H_0 : \beta_2 = 0$ $H_1 : \beta_2 \neq 0$
2. Specify the test statistic and its distribution if the null hypothesis is true

$$F = \frac{(SSE_R - SSE_U)/1}{SSE_U/(75-3)} \sim F_{(1, 72)}$$

3. Set and determine the rejection region

Using $\alpha=.05$, the critical value from the F -distribution is $F_c = F_{(.95,1,72)} = 3.97$.

Thus, H_0 is rejected if $F \geq 3.97$.

6.1 The *F*-Test

4. Calculate the sample value of the test statistic and, if desired, the p -value

$$F = \frac{(SSE_R - SSE_U)/J}{SSE_U/(N - K)} = \frac{(2961.827 - 1718.943)/1}{1718.943/(75 - 3)} = 52.06$$

$$p = P[F_{(1,72)} \geq 52.06] = .0000$$

5. State your conclusion

Since $F = 52.06 \geq F_c$, we reject the null hypothesis and conclude that price does have a significant effect on sales revenue. Alternatively, we reject H_0 because $p = .0000 \leq .05$.

6.1.1 The Relationship Between t - and F -Tests

The elements of an F -test :

1. The null hypothesis consists of one or more equality restrictions J . The null hypothesis may not include any ‘greater than or equal to’ or ‘less than or equal to’ hypotheses.
2. The alternative hypothesis states that one or more of the equalities in the null hypothesis is not true. The alternative hypothesis may not include any ‘greater than’ or ‘less than’ options.
3. The test statistic is the F -statistic
$$F = \frac{(SSE_R - SSE_U)/J}{SSE_U/(N - K)}$$

6.1.1 The Relationship Between t - and F -Tests

4. If the null hypothesis is true, F has the F -distribution with J numerator degrees of freedom and $N-K$ denominator degrees of freedom. The null hypothesis is *rejected* if $F > F_c$, where $F_c = F_{(1-\alpha, J, N-K)}$.
5. When testing a single equality null hypothesis it is perfectly correct to use either the t - or F -test procedure; they are equivalent.

6.2 Testing the Significance of the Model

$$y_i = \beta_1 + x_{i2}\beta_2 + x_{i3}\beta_3 + \cdots + x_{iK}\beta_K + e_i \quad (6.4)$$

$$H_0 : \beta_2 = 0, \beta_3 = 0, \cdots, \beta_K = 0$$

$$H_1 : \text{At least one of the } \beta_k \text{ is nonzero} \quad (6.5)$$

for $k = 2, 3, \dots, K$

$$y_i = \beta_1 + e_i \quad (6.6)$$

6.2 Testing the Significance of the Model

$$SSE_R = \sum_{i=1}^N (y_i - b_1^*)^2 = \sum_{i=1}^N (y_i - \bar{y})^2 = SST$$

$$F = \frac{(SST - SSE) / (K - 1)}{SSE / (N - K)}$$

(6.7)

6.2 Testing the Significance of the Model

Example: Big Andy's sales revenue $S_i = \beta_1 + \beta_2 P_i + \beta_3 A_i + e_i$

1. $H_0 : \beta_2 = 0, \beta_3 = 0$
 $H_1 : \beta_2 \neq 0$ or $\beta_3 \neq 0$, or both are nonzero
2. If the null is true $F = \frac{(SST - SSE) / (3 - 1)}{SSE / (75 - 3)} \square F_{(2,72)}$
3. H_0 is rejected if $F \geq 3.12$
4. $F = \frac{(SST - SSE) / (K - 1)}{SSE / (N - K)} = \frac{(3115.485 - 1718.943) / 2}{1718.943 / (75 - 3)} = 29.25$
5. Since $29.95 > 3.12$ we reject the null and conclude that price or advertising expenditure or both have an influence on sales.

6.3 An Extended Model

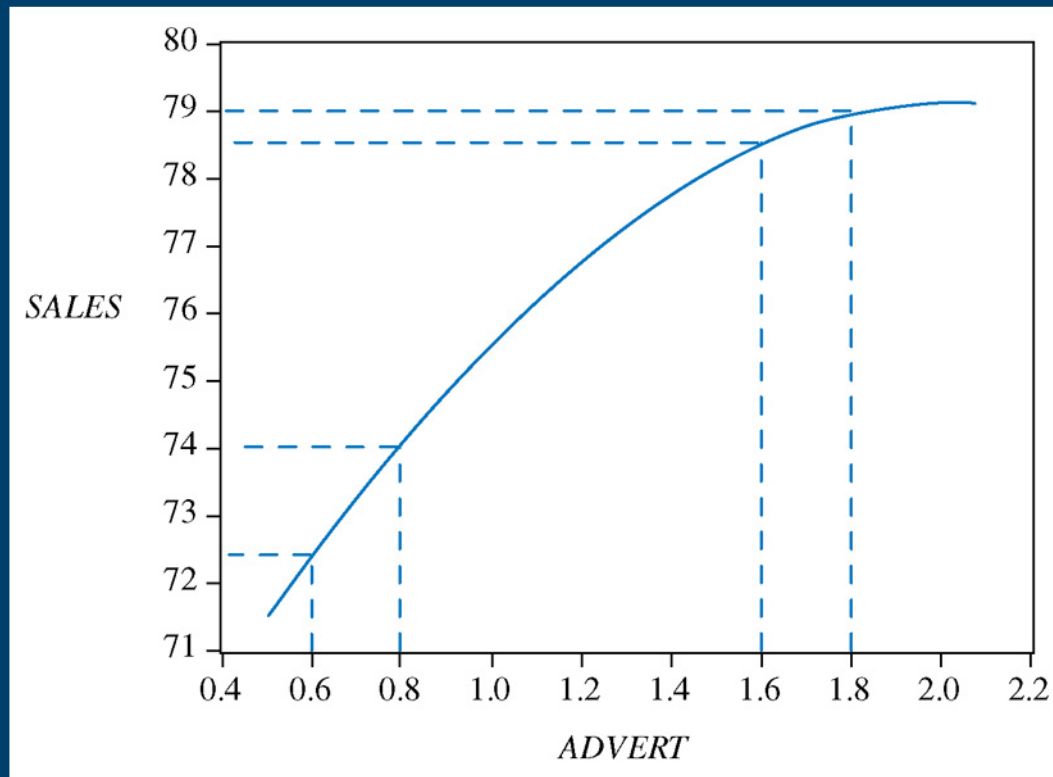


Figure 6.1 A Model Where Sales Exhibits Diminishing Returns to Advertising Expenditure

6.3 An Extended Model

$$S = \beta_1 + \beta_2 P + \beta_3 A + e \quad (6.8)$$

$$S = \beta_1 + \beta_2 P + \beta_3 A + \beta_4 A^2 + e \quad (6.9)$$

$$\frac{\Delta E(S)}{\Delta A} \quad (P \text{ held constant}) = \frac{\partial E(S)}{\partial A} = \beta_3 + 2\beta_4 A \quad (6.10)$$

6.3 An Extended Model

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + e_i$$

$$y_i = S_i, \quad x_{i2} = P_i, \quad x_{i3} = A_i, \quad x_{i4} = A_i^2$$

$$\hat{S}_i = 109.72 - 7.640 P_i + 12.151 A_i - 2.768 A_i^2$$

(se) (6.80) (1.046) (3.556) (.941) (6.11)

6.4 Testing Some Economic Hypotheses

■ 6.4.1 The Significance of Advertising

$$S_i = \beta_1 + \beta_2 P_i + \beta_3 A_i + \beta_4 A_i^2 + e_i \quad (6.12)$$

$$S_i = \beta_1 + \beta_2 P_i + e_i \quad (6.13)$$

$$H_0 : \beta_3 = 0, \beta_4 = 0$$

$$H_1 : \beta_3 \neq 0 \text{ or } \beta_4 \neq 0$$

6.4 Testing Some Economic Hypotheses

$$F = \frac{(SSE_R - SSE_U)/2}{SSE_U/(75 - 4)} \sim F_{(2,71)}$$

$$F_c = F_{(.95,2,71)} = 3.126$$

$$P[F_{(2,71)} > 8.44] = .0005$$

Since $F = 8.44 > F_c = 3.126$, we reject the null hypothesis and conclude that advertising does have a significant effect upon sales revenue.

6.4.2 The Optimal Level of Advertising

Economic theory tells us that we should undertake all those actions for which the marginal benefit is greater than the marginal cost. This optimizing principle applies to Big Andy's Burger Barn as it attempts to choose the optimal level of advertising expenditure.

$$\frac{\Delta E(S)}{\Delta A} \quad (P \text{ held constant}) = \beta_3 + 2\beta_4 A$$

$$\beta_3 + 2\beta_4 A_o = 1$$

$$12.1512 + 2 \times (-2.76796) \hat{A}_o = 1$$

6.4.2 The Optimal Level of Advertising

Big Andy has been spending \$1,900 per month on advertising. He wants to know whether this amount could be optimal.

The null and alternative hypotheses for this test are

$$H_0 : \beta_3 + 2 \times \beta_4 \times 1.9 = 1$$

$$H_1 : \beta_3 + 2 \times \beta_4 \times 1.9 \neq 1$$

$$H_0 : \beta_3 + 3.8 \beta_4 = 1$$

$$H_1 : \beta_3 + 3.8 \beta_4 \neq 1$$

6.4.2 The Optimal Level of Advertising

$$t = \frac{(b_3 + 3.8b_4) - 1}{\text{se}(b_3 + 3.8b_4)}$$

$$\begin{aligned}\text{var}(b_3 + 3.8b_4) &= \text{var}(b_3) + 3.8^2 \times \text{var}(b_4) + 2 \times 3.8 \times \text{cov}(b_3, b_4) \\ &= 12.6463 + 3.8^2 \times .884774 - 2 \times 3.8 \times 3.288746 \\ &= .427967\end{aligned}$$

6.4.2 The Optimal Level of Advertising

$$t = \frac{1.6330 - 1}{\sqrt{.427967}} = \frac{.633}{.65419} = .9676$$

Because $-1.994 < .9676 < 1.994$, we cannot reject the null hypothesis that the optimal level of advertising is \$1,900 per month. There is insufficient evidence to suggest Andy should change his advertising strategy.

6.4.2 The Optimal Level of Advertising

$$S_i = \beta_1 + \beta_2 P_i + (1 - 3.8\beta_4)A_i + \beta_4 A_i^2 + e_i$$

$$(S_i - A_i) = \beta_1 + \beta_2 P_i + \beta_4 (A_i^2 - 3.8A_i) + e_i$$

$$F = \frac{(1552.286 - 1532.084) / 1}{1532.084 / 71} = .9362$$

$$F_c = 3.976 = t_c^2 = (1.994)^2$$

$$\begin{aligned} p\text{-value} &= P[F_{(1,71)} > .9362] \\ &= P[t_{(71)} > .9676] + P[t_{(71)} < -.9676] = .3365 \end{aligned}$$

6.4.2a A One-Tailed Test with More than One Parameter

$$H_0 : \beta_3 + 3.8\beta_4 \leq 1$$

$$H_1 : \beta_3 + 3.8\beta_4 > 1$$

Reject H_0 if $t \geq 1.667$.

$$t = .9676$$

Because $.9676 < 1.667$, we do not reject H_0 .

There is not enough evidence in the data to suggest the optimal level of advertising expenditure is greater than \$1900.

6.4.2 Using Computer Software

$$\begin{aligned}E(S) &= \beta_1 + \beta_2 P + \beta_3 A + \beta_4 A^2 \\ &= \beta_1 + 6\beta_2 + 1.9\beta_3 + 1.9^2 \beta_4 \\ &= 80\end{aligned}$$

$$\begin{aligned}H_0 : \quad &\beta_3 + 3.8\beta_4 = 1 \quad \text{and} \\ &\beta_1 + 6\beta_2 + 1.9\beta_3 + 3.61\beta_4 = 80\end{aligned}$$

$$F = 5.74 \text{ and the } p\text{-value} = .0049$$

6.5 The Use of Nonsample Information

$$\ln(Q) = \beta_1 + \beta_2 \ln(PB) + \beta_3 \ln(PL) + \beta_4 \ln(PR) + \beta_5 \ln(I) \quad (6.14)$$

$$\begin{aligned} \ln(Q) &= \beta_1 + \beta_2 \ln(\lambda PB) + \beta_3 \ln(\lambda PL) + \beta_4 \ln(\lambda PR) + \beta_5 \ln(\lambda I) \\ &= \beta_1 + \beta_2 \ln(PB) + \beta_3 \ln(PL) + \beta_4 \ln(PR) + \beta_5 \ln(I) \\ &\quad + (\beta_2 + \beta_3 + \beta_4 + \beta_5) \ln(\lambda) \end{aligned} \quad (6.15)$$

6.5 The Use of Nonsample Information

$$\beta_2 + \beta_3 + \beta_4 + \beta_5 = 0 \quad (6.16)$$

$$\ln(Q_t) = \beta_1 + \beta_2 \ln(PB_t) + \beta_3 \ln(PL_t) + \beta_4 \ln(PR_t) + \beta_5 \ln(I_t) + e_t \quad (6.17)$$

6.5 The Use of Nonsample Information

Table 6.1 Summary Statistics for Data Used to Estimate Beer Demand

	<i>Q</i>	<i>PB</i>	<i>PL</i>	<i>PR</i>	<i>I</i>
Sample mean	56.11	3.08	8.37	1.25	32,602
Median	54.90	3.11	8.39	1.18	32,457
Maximum	81.70	4.07	9.52	1.73	41,593
Minimum	44.30	1.78	6.95	0.67	25,088
Std. Dev.	7.8574	0.6422	0.7696	0.2983	4,542

6.5 The Use of Nonsample Information

$$\beta_4 = -\beta_2 - \beta_3 - \beta_5$$

$$\ln(Q_t) = \beta_1 + \beta_2 \ln(PB_t) + \beta_3 \ln(PL_t)$$

$$+ (-\beta_2 - \beta_3 - \beta_5) \ln(PR_t) + \beta_5 \ln(I_t) + e_t$$

$$= \beta_1 + \beta_2 (\ln(PB_t) - \ln(PR_t))$$

$$+ \beta_3 (\ln(PL_t) - \ln(PR_t)) + \beta_5 (\ln(I_t) - \ln(PR_t)) + e_t$$

$$= \beta_1 + \beta_2 \ln\left(\frac{PB_t}{PR_t}\right) + \beta_3 \ln\left(\frac{PL_t}{PR_t}\right) + \beta_5 \ln\left(\frac{I_t}{PR_t}\right) + e_t$$

(6.18)

6.5 The Use of Nonsample Information

$$\ln(Q_t) = -4.798 - 1.2994 \ln\left(\frac{PB_t}{PR_t}\right) + 0.1868 \ln\left(\frac{PL_t}{PR_t}\right) + 0.9458 \ln\left(\frac{I_t}{PR_t}\right) \quad (6.19)$$

(se) (0.166) (0.284) (0.427)

$$\begin{aligned} b_4^* &= -b_2^* - b_3^* - b_5^* \\ &= -(-1.2994) - 0.1868 - 0.9458 \\ &= 0.1668 \end{aligned}$$

6.6 Model Specification

■ 6.6.1 Omitted Variables

$$\begin{array}{rcccc} \hat{FAMINC}_i & = & -5534 & + & 3132 HEDU_i & + & 4523 WEDU_i & & \\ (se) & & (11230) & & (803) & & (1066) & & (6.20) \\ (p\text{-value}) & & (.622) & & (.000) & & (.000) & & \end{array}$$

6.6.1 Omitted Variables

Table 6.2 Summary Statistics for Data Used for Family Income Example

	<i>FAMINC</i>	<i>HEDU</i>	<i>WEDU</i>	<i>KL6</i>	X_5	X_6
Sample mean	91213	12.61	12.65	0.14	12.57	25.13
Median	83013	12	12	0	12.60	24.91
Maximum	344146	17	17	2	20.82	37.68
Minimum	9072	4	5	0	2.26	9.37
Std. Dev.	44147	3.035	2.285	0.392	3.427	5.052
Correlation matrix						
<i>FAMINC</i>	1.000					
<i>HEDU</i>	0.355	1.000				
<i>WEDU</i>	0.362	0.594	1.000			
<i>KL6</i>	-0.072	0.105	0.129	1.000		
X_5	0.290	0.836	0.518	0.149	1.000	
X_6	0.351	0.821	0.799	0.160	0.900	1.000

6.6.1 Omitted Variables

$$\begin{array}{rcc} \hat{FAMINC}_i = -26191 + 5155 HEDU_i & & \\ \text{(se)} & (8541) & (658) & (6.21) \\ \text{(p-value)} & (.002) & (.000) & \end{array}$$

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + e_i \quad (6.22)$$

6.6.1 Omitted Variables

$$\text{bias}(b_2^*) = E(b_2^*) - \beta_2 = \beta_3 \frac{\text{cov}(x_2, x_3)}{\text{var}(x_2)} \quad (6.23)$$

	\overline{FAMINC}_i	$= -7755 + 3211HEDU_i + 4777WEDU_i - 14311KL6_i$			
(se)	(11163)	(797)	(1061)	(5004)	(6.24)
(p-value)	(.488)	(.000)	(.000)	(.004)	

6.6.2 Irrelevant Variables

$$\begin{array}{l} \boxed{FAMINC}_i = -7759 + 3340 HEDU_i + 5869 WEDU_i - 14200 KL6_i + 889 X_{i5} - 1067 X_{i6} \\ \text{(se)} \quad (11195) \quad (1250) \quad (2278) \quad (5044) \quad (2242) \quad (1982) \\ \text{(p-value)} \quad (.500) \quad (.008) \quad (.010) \quad (.005) \quad (.692) \quad (.591) \end{array}$$

6.6.3 Choosing the Model

1. Choose variables and a functional form on the basis of your theoretical and general understanding of the relationship.

where $\beta_2^* = c\beta_2$ and $x^* = x/c$

2. If an estimated equation has coefficients with unexpected signs, or unrealistic magnitudes, they could be caused by a misspecification such as the omission of an important variable.

6.6.3 Choosing the Model

3. One method for assessing whether a variable or a group of variables should be included in an equation is to perform significance tests. That is, t -tests for hypotheses such as $H_0 : \beta_3 = 0$ or F -tests for hypotheses such as $H_0 : \beta_3 = \beta_4 = 0$. Failure to reject hypotheses such as these can be an indication that the variable(s) are irrelevant.
4. The adequacy of a model can be tested using a general specification test known as RESET.

6.6.3a The RESET Test

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + e_i$$

$$\hat{y}_i = b_1 + b_2 x_{i2} + b_3 x_{i3}$$

(6.25)

6.6.3a The RESET Test

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \gamma_1 \hat{y}_i^2 + e_i \quad (6.26)$$

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \gamma_1 \frac{y_i}{T} + \gamma_2 y_i^3 + e_i \quad (6.27)$$

6.6.3a *The RESET Test*

$$H_0 : \gamma_1 = 0 \qquad F = 5.984 \qquad p\text{-value} = .015$$

$$H_0 : \gamma_1 = \gamma_2 = 0 \qquad F = 3.123 \qquad p\text{-value} = .045$$

6.7 Poor data, Collinearity and Insignificance

■ 6.7.1 The Consequences of Collinearity

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + e_i$$

$$\text{var}(b_2) = \frac{\sigma^2}{(1 - r_{23}^2) \sum_{i=1}^N (x_{i2} - \bar{x}_2)^2} \quad (6.28)$$

6.7.1 The Consequences of Collinearity

The effects of imprecise information:

1. When estimator standard errors are large, it is likely that the usual t -tests will lead to the conclusion that parameter estimates are not significantly different from zero. This outcome occurs despite possibly high R^2 or F -values indicating significant explanatory power of the model as a whole.

6.7.1 The Consequences of Collinearity

2. The estimators may be very sensitive to the addition or deletion of a few observations, or the deletion of an apparently insignificant variable.
3. Despite the difficulties in isolating the effects of individual variables from such a sample, accurate forecasts may still be possible if the nature of the collinear relationship remains the same within the new (future) sample observations.

6.7.2 An Example

MPG = miles per gallon

CYL = number of cylinders

ENG = engine displacement in cubic inches

WGT = vehicle weight in pounds

$$\begin{array}{r} \square \\ \text{(se)} \\ \text{(} p\text{-value)} \end{array} \begin{array}{l} \text{MPG}_i = 42.9 - 3.558 \text{CYL}_i \\ \\ \\ \end{array} \begin{array}{l} \\ \\ \\ \end{array} \begin{array}{l} \\ \\ \\ \end{array}$$

(se)	(.83)	(.146)	
(<i>p</i> -value)	(.000)	(.000)	

6.7.2 An Example

$$\begin{array}{rcccc} \square & \text{MPG}_i & = & 44.4 & - & .268 & \text{CYL}_i & - & .0127 & \text{ENG}_i & - & .00571 & \text{WGT}_i \\ & (\text{se}) & & (1.5) & & (.413) & & & (.0083) & & & (.0071) \\ & (p\text{-value}) & & (.000) & & (.517) & & & (.125) & & & (.000) \end{array}$$

6.7.3 Identifying and Mitigating Collinearity

Identifying Collinearity

1. Examining pairwise correlations.
2. Using auxiliary regression

$$x_{i2} = a_1 x_{i1} + a_3 x_{i3} + \cdots + a_K x_{iK} + error$$

If the R^2 from this artificial model is high, above .80 say, the implication is that a large portion of the variation in x_{i2} is explained by variation in the other explanatory variables.

6.7.3 Identifying and Mitigating Collinearity

Mitigating Collinearity

1. Obtain more information and include it in the analysis.
2. Introduce *nonsample* information in the form of restrictions on the parameters.

6.8 Prediction

$$y_i = \beta_1 + x_{i2}\beta_2 + x_{i3}\beta_3 + e_i \quad (6.29)$$

$$y_0 = \beta_1 + x_{02}\beta_2 + x_{03}\beta_3 + e_0$$

$$\hat{y}_0 = b_1 + x_{02}b_2 + x_{03}b_3$$

6.8 Prediction

$$\begin{aligned}\text{var}(f) &= \text{var}[(\beta_1 + \beta_2 x_{02} + \beta_3 x_{03} + e_0) - (b_1 + b_2 x_{02} + b_3 x_{03})] \\ &= \text{var}(e_0 - b_1 - b_2 x_{02} - b_3 x_{03}) \\ &= \text{var}(e_0) + \text{var}(b_1) + x_{02}^2 \text{var}(b_2) + x_{03}^2 \text{var}(b_3) \\ &\quad + 2x_{02} \text{cov}(b_1, b_2) + 2x_{03} \text{cov}(b_1, b_3) + 2x_{02}x_{03} \text{cov}(b_2, b_3)\end{aligned}$$

Keywords

- a single null hypothesis with more than one parameter
- auxiliary regressions
- collinearity
- F -test
- irrelevant variable
- nonsample information
- omitted variable
- omitted variable bias
- overall significance of a regression model
- regression specification error test (RESET)
- restricted least squares
- restricted sum of squared errors
- single and joint null hypotheses
- unrestricted sum of squared errors

Chapter 6 Appendices

- Appendix 6A Chi-Square and F -tests: More Details
- Appendix 6B Omitted Variable Bias: A Proof

Appendix 6A

Chi-Square and F -tests: More Details

$$F = \frac{(SSE_R - SSE_U) / J}{SSE_U / (N - K)} \quad (6A.1)$$

$$V_1 = \frac{(SSE_R - SSE_U)}{\sigma^2} \square \chi_{(J)}^2 \quad (6A.2)$$

$$\hat{V}_1 = \frac{(SSE_R - SSE_U)}{\hat{\sigma}^2} \square \chi_{(J)}^2 \quad (6A.3)$$

$$V_2 = \frac{(N - K)\hat{\sigma}^2}{\sigma^2} \square \chi_{(N-K)}^2 \quad (6A.4)$$

Appendix 6A

Chi-Square and F -tests: More Details

$$F = \frac{V_1 / m_1}{V_2 / m_2} \square F(m_1, m_2)$$

$$\frac{\frac{(SSE_R - SSE_U) / J}{\sigma^2}}{\frac{(N - K) \hat{\sigma}^2 / (N - K)}{\sigma^2}} = \frac{[SSE_R - SSE_U] / J}{\hat{\sigma}^2} \square F_{(J, N-K)} \quad (6A.5)$$

Appendix 6A

Chi-Square and F -tests: More Details

$$F = \frac{\hat{V}_1}{J}$$

$$H_0 : \beta_3 = \beta_4 = 0$$

$$S_i = \beta_1 + \beta_2 P_i + \beta_3 A_i + \beta_4 A_i^2 + e_i$$

$$F = 8.44 \quad p\text{-value} = .0005$$

$$\chi^2 = 16.88 \quad p\text{-value} = .0002$$

Appendix 6A

Chi-Square and F -tests: More Details

$$H_0 : \beta_3 + 3.8\beta_4 = 1$$

$$F = .936$$

$$p\text{-value} = .3365$$

$$\chi^2 = .936$$

$$p\text{-value} = .3333$$

Appendix 6B

Omitted Variable Bias: A Proof

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + e_i$$

$$y_i = \beta_1 + \beta_2 x_{i2} + v_i$$

Appendix 6B

Omitted Variable Bias: A Proof

$$b_2^* = \frac{\sum (x_{i2} - \bar{x}_2)(y_i - \bar{y})}{\sum (x_{i2} - \bar{x}_2)^2} = \beta_2 + \sum w_i v_i \quad (6B.1)$$

$$w_i = \frac{(x_{i2} - \bar{x}_2)}{\sum (x_{i2} - \bar{x}_2)^2}$$

$$b_2^* = \beta_2 + \beta_3 \sum w_i x_{i3} + \sum w_i e_i$$

Appendix 6B

Omitted Variable Bias: A Proof

$$\begin{aligned} E(b_2^*) &= \beta_2 + \beta_3 \sum w_i x_{i3} \\ &= \beta_2 + \beta_3 \frac{\sum (x_{i2} - \bar{x}_2) x_{i3}}{\sum (x_{i2} - \bar{x}_2)^2} \\ &= \beta_2 + \beta_3 \frac{\sum (x_{i2} - \bar{x}_2)(x_{i3} - \bar{x}_3)}{\sum (x_{i2} - \bar{x}_2)^2} \\ &= \beta_2 + \beta_3 \frac{\overline{\text{cov}}(x_2, x_3)}{\overline{\text{var}}(x_2)} \neq \beta_2 \end{aligned}$$