Chapter 6 Further Inference in the Multiple Regression Model

Walter R. Paczkowski Rutgers University

Principles of Econometrics, 4t^h Edition Chapter 6: Further Inference in the Multiple Regression Model

Chapter Contents

- 6.1 Joint Hypothesis Testing
- 6.2 The Use of Nonsample Information
- 6.3 Model Specification
- 6.4 Poor Data, Collinearity, and Insignificance
- 6.5 Prediction

Economists develop and evaluate theories about economic behavior

- Hypothesis testing procedures are used to test these theories
- The theories that economists develop sometimes provide nonsample information that can be used along with the information in a sample of data to estimate the parameters of a regression model
- A procedure that combines these two types of information is called restricted least squares

6.1 Testing Joint Hypotheses

A null hypothesis with multiple conjectures, expressed with more than one equal sign, is called a joint hypothesis

- 1. Example: Should a group of explanatory variables should be included in a particular model?
- 2. Example: Does the quantity demanded of a product depend on the prices of substitute goods, or only on its own price?

> 6.1 Testing Joint Hypotheses

Eq. 6.1

Both examples are of the form:

$$H_0: \beta_4 = 0, \beta_5 = 0, \beta_6 = 0$$

- The joint null hypothesis in Eq. 6.1 contains three conjectures (three equal signs): $\beta_4 = 0$, $\beta_5 = 0$, and $\beta_6 = 0$
- A test of H_0 is a joint test for whether all three conjectures hold simultaneously

6.1.1 Testing the Effect of Advertising: The *F*-Test

6.1 Joint Hypothesis T<u>esting</u>

Eq. 6.2

- $SALES = \beta_1 + \beta_2 PRICE + \beta_3 ADVERT + \beta_4 ADVERT^2 + e$
 - Test whether or not advertising has an effect on sales – but advertising is in the model as two variables

Consider the model:

6.1.1 Testing the Effect of Advertising: The *F*-Test

- Advertising will have no effect on sales if $\beta_3 = 0$ and $\beta_4 = 0$
- Advertising will have an effect if $\beta_3 \neq 0$ or $\beta_4 \neq 0$ or if both β_3 and β_4 are nonzero
- The null hypotheses are:

 $H_0: \beta_3 = 0, \beta_4 = 0$ $H_1: \beta_3 \neq 0 \text{ or } \beta_4 \neq 0 \text{ or both are nonzero}$

6.1.1 Testing the Effect of Advertising: The *F*-Test

- Relative to the null hypothesis H_0 : $\beta_3 = 0$, $\beta_4 = 0$ the model in Eq. 6.2 is called the **unrestricted model**
 - The restrictions in the null hypothesis have not been imposed on the model
 - It contrasts with the restricted model, which is obtained by assuming the parameter restrictions in H_0 are true

6.1.1 Testing the Effect of Advertising: The *F*-Test

Eq. 6.3

When H_0 is true, $\beta 3 = 0$ and $\beta_4 = 0$, and ADVERTand $ADVERT^2$ drop out of the model

 $SALES = \beta_1 + \beta_2 PRICE + e$

The *F*-test for the hypothesis H₀: β₃ = 0, β₄ = 0 is based on a comparison of the sums of squared errors (sums of squared least squares residuals) from the unrestricted model in Eq. 6.2 and the restricted model in Eq. 6.3

- Shorthand notation for these two quantities is SSE_U and SSE_R , respectively

6.1.1 Testing the Effect of Advertising: The *F*-Test

The F-statistic determines what constitutes a large reduction or a small reduction in the sum of squared errors

Eq. 6.4

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

where *J* is the number of restrictions, *N* is the number of observations and *K* is the number of coefficients in the unrestricted model

6.1.1 Testing the Effect of Advertising: The *F*-Test

- If the null hypothesis is true, then the statistic F has what is called an F-distribution with J numerator degrees of freedom and N K denominator degrees of freedom
- If the null hypothesis is not true, then the difference between SSE_R and SSE_U becomes large
 - The restrictions placed on the model by the null hypothesis significantly reduce the ability of the model to fit the data

6.1.1 Testing the Effect of Advertising: The *F*-Test

The *F*-test for our sales problem is:

- 1. Specify the null and alternative hypotheses:
 - The joint null hypothesis is H₀: β₃ = 0, β₄ = 0. The alternative hypothesis is H₀: β₃ ≠ 0 or β₄ ≠ 0 both are nonzero
- 2. Specify the test statistic and its distribution if the null hypothesis is true:
 - Having two restrictions in H_0 means J = 2

• Also, recall that
$$N = 75$$
:

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

6.1.1 Testing the Effect of Advertising: The *F*-Test

■ The *F*-test for our sales problem is (Continued):

- 3. Set the significance level and determine the rejection region
- 4. Calculate the sample value of the test statistic and, if desired, the *p*-value

$$F = \frac{\left(SSE_R - SSE_U\right)/J}{SSE_U/(N-K)} = \frac{\left(1896.391 - 1532.084\right)/2}{1532.084/(75-4)} = 8.44$$

• The corresponding *p*-value is $p = P(F_{(2,71)} > 8.44) = 0.0005$

6.1.1 Testing the Effect of Advertising: The *F*-Test

■ The *F*-test for our sales problem is (Continued):

- 5. State your conclusion
 - Since $F = 8.44 > F_c = 3.126$, we reject the null hypothesis that both $\beta_3 = 0$ and $\beta_4 = 0$, and conclude that at least one of them is not zero
 - -Advertising does have a significant effect upon sales revenue

6.1.2 Testing the Significance of the Model

Eq. 6.5

Consider again the general multiple regression model with (K - 1) explanatory variables and K unknown coefficients

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + L + \beta_K x_K + e$$

To examine whether we have a viable explanatory model, we set up the following null and alternative hypotheses:

Eq. 6.6

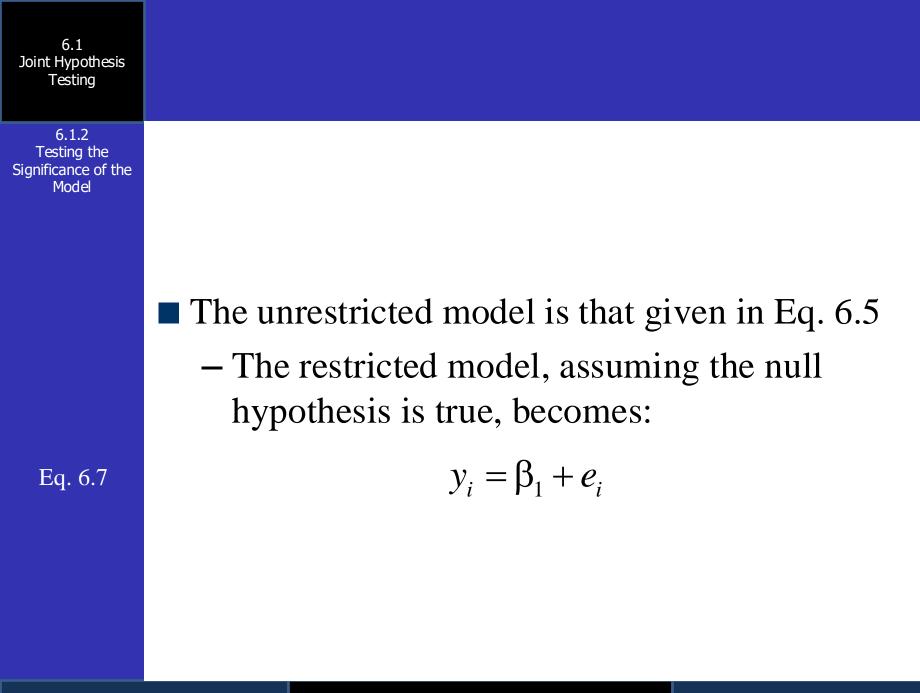
$$H_0: \beta_2 = 0, \beta_3 = 0, K, \beta_K = 0$$

 H_1 : At least one of the β_k is nonzero for k = 2,3,K K

6.1.2 Testing the Significance of the Model

Since we are testing whether or not we have a viable explanatory model, the test for Eq. 6.6 is sometimes referred to as a **test of the overall significance of the regression model**.

- Given that the *t*-distribution can only be used to test a single null hypothesis, we use the *F*-test for testing the joint null hypothesis in Eq. 6.6



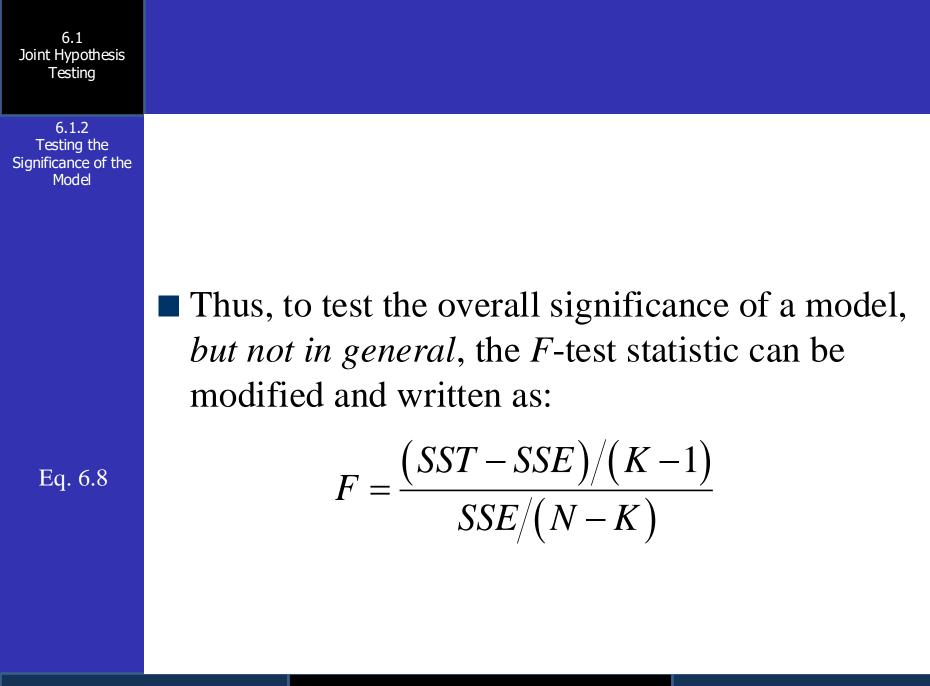
6.1.2 Testing the Significance of the Model

The least squares estimator of β_1 in this restricted model is:

$$b_1^* = \sum_{i=1}^N y_i / N = \overline{y}$$

The restricted sum of squared errors from the hypothesis Eq. 6.6 is:

$$SSE_{R} = \sum_{i=1}^{N} (y_{i} - b_{1}^{*})^{2} = \sum_{i=1}^{N} (y_{i} - \overline{y})^{2} = SST$$



6.1.2 Testing the Significance of the Model

For our problem, note:

1. We are testing:

 $H_0: \beta_2 = 0, \beta_3 = 0, \beta_4 = 0$ $H_1: At \ least \ one \ of \ \beta_2 \ or \ \beta_3 \ or \ \beta_4 \ is \ nonzero$

2. If H_0 is true:

$$F = \frac{(SST - SSE)/(4 - 1)}{SSE/(75 - 4)} \sim F_{(3,71)}$$

6.1.2 Testing the Significance of the Model

■ For our problem, note (Continued):

- 3. Using a 5% significance level, we find the critical value for the *F*-statistic with (3,71) degrees of freedom is $F_c = 2.734$.
 - Thus, we reject H_0 if F ≥ 2.734 .
- 4. The required sums of squares are SST = 3115.482 and SSE = 1532.084 which give an *F*-value of:

$$F = \frac{\left(SST - SSE\right) / (K - 1)}{SSE / (N - K)} = \frac{\left(3115.482 - 1532.084\right) / 3}{1532.084 / (75 - 4)} = 24.459$$

• *p*-value =
$$P(F \ge 24.459) = 0.0000$$

6.1.2 Testing the Significance of the Model

■ For our problem, note (Continued):

- 5. Since 24.459 > 2.734, we reject H_0 and conclude that the estimated relationship is a significant one
 - Note that this conclusion is consistent with conclusions that would be reached using separate *t*-tests for the significance of each of the coefficients

6.1.3 Relationship Between the *t*- and *F*-Tests

We used the *F*-test to test whether $\beta_3 = 0$ and $\beta_4 = 0$ in:

Eq. 6.9

 $SALES = \beta_1 + \beta_2 PRICE + \beta_3 ADVERT + \beta_4 ADVERT^2 + e$

- Suppose we want to test if *PRICE* affects SALES $H_0: \beta_1 = 0$

$$H_1: \beta_1 \neq 0$$

Eq. 6.10

 $SALES = \beta_1 + \beta_3 ADVERT + \beta_4 ADVERT^2 + e$

6.1.3 Relationship Between the *t*- and *F*-Tests

■ The *F*-value for the restricted model is:

$$F = \frac{\left(SSE_R - SSE_U\right)/J}{SSE_U/(N-K)} = \frac{\left(2683.411 - 1532.084\right)/1}{1532.084/(75-4)} = 53.355$$

- The 5% critical value is $F_c = {}_{F(0.95, 1, 71)} = 3.976$ - We reject H_0 : $\beta_2 = 0$

6.1.3 Relationship Between the *t*- and *F*-Tests

Using the *t*-test:

 $SALES = 109.72 - 7.640 PRICE + 12.151 ADVERT - 2.768 ADVERT^{2}$ (se) (6.80) (1.046) (3.556) (0.941)

- The *t*-value for testing H_0 : $\beta_2 = 0$ against H_1 : $\beta_2 \neq 0$ is t = 7.640/1.045939 = 7.30444
- Its square is $t = (7.30444)^2 = 53.355$, identical to the *F*-value

6.1.3 Relationship Between the *t*- and *F*-Tests

■ The elements of an *F*-test

- 1. The null hypothesis H_0 consists of one or more equality restrictions on the model parameters β_k
- 2. The alternative hypothesis states that one or more of the equalities in the null hypothesis is not true
- 3. The test statistic is the F-statistic in (6.4)
- 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
- 5. When testing a single equality null hypothesis, it is perfectly correct to use either the *t* or *F*-test procedure: they are equivalent

6.1.4 More General *F*-Tests

> The conjectures made in the null hypothesis were that particular coefficients are equal to zero

- The *F*-test can also be used for much more general hypotheses
- Any number of conjectures ($\leq K$) involving linear hypotheses with equal signs can be tested

Consider the issue of testing:
$\beta_3 + 2\beta_4 ADVERT_0 = 1$
- If $ADVERT_0 = $1,900$ per month, then:
$H_0: \beta_3 + 2 \times \beta_4 \times 1.9 = 1 \qquad H_1: \beta_3 + 2 \times \beta_4 \times 1.9 \neq 0$ or
$H_0: \beta_3 + 3.8\beta_4 = 1$ $H_1: \beta_3 + 3.8\beta_4 \neq 1$

1

6.1.4 More General *F*-Tests

Eq. 6.12

Note that when H_0 is true, $\beta_3 = 1 - 3.8\beta_4$ so that: $SALES = \beta_1 + \beta_2 PRICE + (1 - 3.8\beta_4) ADVERT + \beta_4 ADVERT^2 + e$ or $(SALES - ADVERT) = \beta_1 + \beta_2 PRICE + \beta_4 (ADVERT^2 - 3.8ADVERT) + e$

6.1.4 More General *F*-Tests

■ The calculated value of the *F*-statistic is:

$$F = \frac{\left(1552.286 - 1532.084\right)/1}{1532.084/71} = 0.9362$$

- For $\alpha = 0.05$, the critical value is $F_c = 3.976$ Since $F = 0.9362 < F_c = 3.976$, we do not reject H_0
- We conclude that an advertising expenditure of \$1,900 per month is optimal is compatible with the data

6 1

The *t*-value is t = 0.9676-F = 0.9362 is equal to $t^2 = (0.9676)^2$ - The *p*-values are identical: $p - value = P(F_{(1,71)} > 0.9362)$ $= P\left(t_{(71)} > 0.9676\right) + P\left(t_{(71)} < -0.9676\right)$ = 0.3365

> 6.1.4a One-tail Test

Eq. 6.13

■ Suppose we have:

 $H_0: \beta_3 + 3.8\beta_4 \le 1$ $H_1: \beta_3 + 3.8\beta_4 > 1$

In this case, we can no longer use the F-test

- Because $F = t^2$, the *F*-test cannot distinguish between the left and right tails as is needed for a one-tail test
- We restrict ourselves to the *t*-distribution when considering alternative hypotheses that have inequality signs such as < or >

6.1.5 Using Computer Software

> Most software packages have commands that will automatically compute *t*- and *F*-values and their corresponding *p*-values when provided with a null hypothesis

These tests belong to a class of tests called
 Wald tests

6.1.5 Using Computer Software

Suppose we conjecture that:

 $E(SALES) = \beta_1 + \beta_2 PRICE + \beta_3 ADVERT + \beta_4 ADVERT^2$ $= \beta_1 + 6\beta_2 + 1.9\beta_3 + 1.9^2\beta_4$ = 80

– We formulate the joint null hypothesis:

 $H_0: \beta_3 + 3.8\beta_4 = 1$ and $\beta_1 + 6\beta_2 + 1.9\beta_3 + 3.61\beta_4 = 80$

- Because there are J = 2 restrictions to test jointly, we use an F-test
 - A *t*-test is not suitable

6.2 The Use of Nonsample Information

Page 36

In many estimation problems we have information over and above the information contained in the sample observations

- This nonsample information may come from many places, such as economic principles or experience
- When it is available, it seems intuitive that we should find a way to use it

Consider the log-log functional form for a demand model for beer:

Eq. 6.14

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

This model is a convenient one because it precludes infeasible negative prices, quantities, and income, and because the coefficients β₂, β₃, β₄, and β₅ are elasticities

- A relevant piece of nonsample information can be derived by noting that if all prices and income go up by the same proportion, we would expect there to be no change in quantity demanded
 - For example, a doubling of all prices and income should not change the quantity of beer consumed
 - This assumption is that economic agents do not suffer from "money illusion"

Eq. 6.15

Having all prices and income change by the same proportion is equivalent to multiplying each price and income by a constant, say λ:

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

To have no change in ln(Q) when all prices and income go up by the same proportion, it must be true that:

Eq. 6.16

$$\beta_2 + \beta_3 + \beta_4 + \beta_5 = 0$$

- We call such a restriction **nonsample information** To estimate a model, start with:

Eq. 6.17

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

- Solve the restriction for one of the parameters, say β_4 :

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

Eq. 6.18

Substituting gives:

$$\ln(Q) = \beta_1 + \beta_2 \ln(PB) + \beta_3 \ln(PL) + (-\beta_2 - \beta_3 - \beta_5) \ln(PR) + \beta_5 \ln(I) + e$$

$$= \beta_1 + \beta_2 \left[\ln(PB) - \ln(PR) \right] + \beta_3 \left[\ln(PL) - \ln(PR) \right]$$

$$+ \beta_5 \left[\ln(I) - \ln(PR) \right] + e$$

$$= \beta_1 + \beta_2 \ln\left(\frac{PB}{PR}\right) + \beta_3 \ln\left(\frac{PL}{PR}\right) + \beta_5 \ln\left(\frac{I}{PR}\right) + e$$

To get least squares estimates that satisfy the parameter restriction, called restricted least squares estimates, we apply the least squares estimation procedure directly to the restricted model:

$$\ln(Q) = -4.798 - 1.2994 \ln\left(\frac{PB}{PR}\right) + 0.1868 \ln\left(\frac{PL}{PR}\right) + 0.9458 \ln\left(\frac{I}{PR}\right)$$
(se) (0.166) (0.284) (0.427)

Eq. 6.19

- Let the restricted least squares estimates in Eq.
 6.19 be denoted by b*₁, b*₂, b*₃, and b*₅
 - To obtain an estimate for β_4 , we use the restriction:

$$b_4^* = -b_2^* - b_3^* - b_5^* = -(-1.2994) - 0.1868 - 0.9458 = 0.1668$$

By using the restriction *within* the model, we have ensured that the estimates obey the constraint:

$$b_2^* + b_3^* + b_4^* + b_4^* = 0$$

Properties of this restricted least squares estimation procedure:

- The restricted least squares estimator is biased, unless the constraints we impose are exactly true
- 2. The restricted least squares estimator is that its variance is smaller than the variance of the least squares estimator, whether the constraints imposed are true or not

In any econometric investigation, choice of the model is one of the first steps

- What are the important considerations when choosing a model?
- What are the consequences of choosing the wrong model?
- Are there ways of assessing whether a model is adequate?

6.3.1 Omitted Variables

> It is possible that a chosen model may have important variables omitted

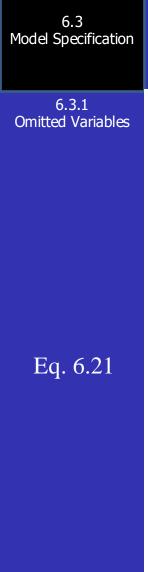
 Our economic principles may have overlooked a variable, or lack of data may lead us to drop a variable even when it is prescribed by economic theory

6.3.1 Omitted Variables

Consider the model:

Eq. 6.20

FAMINC = -5534 + 3132HEDU + 4523WEDU(se)(11230)(803)(1066)(p-value)(0.622)(0.000)(0.000)



If we incorrectly omit wife's education:

FAMINC = -26191 + 5155HEDU(se) (8541) (658) (*p*-value) (0.002)(0.000)

6.3.1 Omitted Variables

> Relative to Eq. 6.20, omitting WEDU leads us to overstate the effect of an extra year of education for the husband by about \$2,000

- Omission of a relevant variable (defined as one whose coefficient is nonzero) leads to an estimator that is biased
- This bias is known as **omitted-variable bias**

6.3.1 Omitted Variables

Eq. 6.22

Write a general model as:

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + e$$

- Omitting x_3 is equivalent to imposing the restriction $\beta_3 = 0$
 - It can be viewed as an example of imposing an incorrect constraint on the parameters

6.3.1 Omitted Variables

The bias is:

Eq. 6.23

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

6.3 Model Specification	Table 6.	1 Correlation 1	Matrix for Va	ariables Used	l in Family	Income Exa	ample
6.3.1 Omitted Variables							
						,	L
		FAMINC	HEDU	WEDU	KL6	X5	х ₆
	FAMINC	<i>FAMINC</i> 1.000	HEDU	WEDU			
	FAMINC HEDU		<i>HEDU</i> 1.000	WEDU			
		1.000		<i>WEDU</i> 1.000			
	HEDU	1.000 0.355	1.000				
	HEDU WEDU	1.000 0.355 0.362	1.000 0.594	1.000	KL6		

6.3.1 **Omitted Variables**

Note that:

- 1. $\beta_3 > 0$ because husband's education has a positive effect on family income.
- 2. $\overline{cov}(x_2, x_3)$ because husband's and wife's levels of education are positively correlated.
- Thus, the bias is positive

6.3.1 Omitted Variables

Eq. 6.24

Now consider the model:

FAMINC = -7755 + 3212HEDU + 4777WEDU - 14311KL6(se)(11163)(797)(1061)(5004)(p-value)(0.488)(0.000)(0.000)(0.004)

- Notice that the coefficient estimates for *HEDU* and *WEDU* have not changed a great deal
 - This outcome occurs because *KL6* is not highly correlated with the education variables

> 6.3.2 Irrelevant Variables

> > You to think that a good strategy is to include as many variables as possible in your model.

 Doing so will not only complicate your model unnecessarily, but may also inflate the variances of your estimates because of the presence of irrelevant variables.

> 6.3.2 Irrelevant Variables

> > You to think that a good strategy is to include as many variables as possible in your model.

 Doing so will not only complicate your model unnecessarily, but may also inflate the variances of your estimates because of the presence of irrelevant variables.

> 6.3.2 Irrelevant Variables

Consider the model:

FAMINC = $-7759 + 3340HEDU + 5869WEDU - 14200KL6 + 889X_5 - 1067X_6$ (se)(11195)(1250)(2278)(5044)(2242)(1982)(p-value)(0.500)(0.008)(0.010)(0.005)(0.692)(0.591)

 The inclusion of irrelevant variables has reduced the precision of the estimated coefficients for other variables in the equation

6.3.3 Choosing the Model

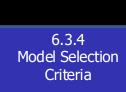
Some points for choosing a model:

- 1. Choose variables and a functional form on the basis of your theoretical and general understanding of the relationship
- 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
- 3. One method for assessing whether a variable or a group of variables should be included in an equation is to perform significance tests

> 6.3.3 Choosing the Model

> > Some points for choosing a model (Continued):

- 4. Consider various model selection criteria
- 5. The adequacy of a model can be tested using a general specification test known as RESET



There are three main model selection criteria:

- 1. R^2
- 2. AIC
- *3. SC* (*BIC*)

6.3.4 Model Selection Criteria

> • A common feature of the criteria we describe is that they are suitable only for comparing models with the same dependent variable, not models with different dependent variables like y and ln(y)

> 6.3.4a The Adjusted Coefficient of Determination

- The problem is that *R*² can be made large by adding more and more variables, even if the variables added have no justification
 - Algebraically, it is a fact that as variables are added the sum of squared errors SSE goes down, and thus R^2 goes up
 - If the model contains N 1 variables, then $R^2 = 1$

> 6.3.4a The Adjusted Coefficient of Determination

> > An alternative measure of <u>goo</u>dness of fit called the adjusted- R^2 , denoted as R^2 :

Eq. 6.25

$$\overline{R^2} = 1 - \frac{SSE/(N-K)}{SST/(N-1)}$$

> 6.3.4b Information Criteria

The Akaike information criterion (AIC) is given by:

Eq. 6.26

$$AIC = \ln\left(\frac{SSE}{N}\right) + \frac{2K}{N}$$

Principles of Econometrics, 4t^h Edition Chapter 6: Further Inference in the Multiple Regression Model

> 6.3.4b Information Criteria

Schwarz criterion (*SC*), also known as the **Bayesian information criterion** (*BIC*) is given by: $SC = \ln\left(\frac{SSE}{N}\right) + \frac{K\ln(N)}{N}$

Eq. 6.27

Page 68

6.3 Model Specification	Table 6.2 Goodness-of-Fit and Information Criteria for Family Income Example							
6.3.4b Information Criteria								
	Included Variables	R^2	\overline{R}^2	AIC	SC			
	HEDU	0.1258	0.1237	21.262	21.281			
	HEDU, WEDU	0.1613	0.1574	21.225	21.253			
	HEDU, WEDU, KL6	0.1771	0.1714	21.211	21.248			
	HEDU, WEDU, KL6, X5, X6	0.1778	0.1681	21.219	21.276			

6.3.5 RESET

- A model could be misspecified if:
 - we have omitted important variables
 - included irrelevant ones
 - chosen a wrong functional form
 - have a model that violates the assumptions of the multiple regression model

> 6.3.5 RESET

RESET (REgression Specification Error Test) is designed to detect omitted variables and incorrect functional form

6.3.5 RESET

6.3 Model Specification

Suppose we have the model:

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + e$$

$$\hat{y} = b_1 + b_2 x_2 + b_3 x_3$$

Eq. 6.28

6.3.5 RESET

6.3 Model Specification

Now consider the following two artificial models:

Eq. 6.29

Eq. 6.30

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \gamma_1 \hat{y}^2 + e$$

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \gamma_1 \hat{y}^2 + \gamma_1 \hat{y}^{3s} + e$$

6.3.5 RESET

6.3 Model Specification

In Eq. 6.29 a test for misspecification is a test of *H*₀:γ₁ = 0 against the alternative *H*₁:γ₁ ≠ 0
 In Eq. 6.30, testing *H*₀:γ₁ = γ₂ = 0 against *H*₁: γ₁ ≠ 0 and/or γ₂ ≠ 0 is a test for misspecification

6.3 Model Specification

> 6.3.5 RESET

> > Applying RESET to our problem (Eq. 6.24), we get:

 $H_0: \gamma_1 = 0 \qquad F = 5.984 \qquad p - value = 0.015$ $H_0: \gamma_1 = \gamma_2 = 0 \quad F = 3.123 \qquad p - value = 0.045$

 In both cases the null hypothesis of no misspecification is rejected at a 5% significance level

When data are the result of an uncontrolled experiment, many of the economic variables may move together in systematic ways

Such variables are said to be collinear, and the problem is labeled collinearity

6.4.1 The Consequences of Collinearity

Consider the model:

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + e$$

– The variance of the least squares estimator for β_2 is:

$$\operatorname{var}(b_{2}) = \frac{\sigma^{2}}{\left(1 - r_{23}^{2}\right) \sum_{i=1}^{n} (x_{2} - \overline{x}_{2})^{2}}$$

Eq. 6.31

6.4.1 The Consequences of Collinearity

> Exact or extreme collinearity exists when x_2 and x_3 are perfectly correlated, in which case $r_{23} = 1$ and var(b₂) goes to infinity

- Similarly, if x_2 exhibits no variation $\sum (x_2 \overline{x}_2)^2$ equals zero and var(b₂) again goes to infinity
 - In this case *x*₂ is collinear with the constant term

6.4.1 The Consequences of Collinearity

> In general, whenever there are one or more exact linear relationships among the explanatory variables, then the condition of exact collinearity exists

- In this case the least squares estimator is not defined
- We cannot obtain estimates of β_k 's using the least squares principle

6.4.1 The Consequences of Collinearity ■ The effects of this imprecise information are:

- 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
- 2. Estimators may be very sensitive to the addition or deletion of a few observations, or to the deletion of an apparently insignificant variable
- 3. Accurate forecasts may still be possible if the nature of the collinear relationship remains the same within the out-of-sample observations

> 6.4.2 An Example

A regression of *MPG* on *CYL* yields:

MPG = 42.9 - 3.558CYL(se) (0.83) (0.146) (p-value)(0.000)(0.000)

– Now add *ENG* and *WGT*:

MPG = 44.4 - 0.268CYL - 0.0127ENG - 0.00571WGT(se) (1.5) (0.413) (0.0083) (0.00071) (p-value)(0.000)(0.517) (0.125) (0.000)

6.4.3 Identifying and Mitigating Collinearity

One simple way to detect collinear relationships is to use sample correlation coefficients between pairs of explanatory variables

- These sample correlations are descriptive measures of linear association
- However, in some cases in which collinear relationships involve more than two of the explanatory variables, the collinearity may not be detected by examining pairwise correlations

6.4.3 Identifying and Mitigating Collinearity

Try an auxiliary model:

$$x_2 = a_1 x_1 + a_3 x_3 + L + a_K x_K + error$$

- If R^2 from this artificial model is high, above 0.80, say, the implication is that a large portion of the variation in x_2 is explained by variation in the other explanatory variables

6.4.3 Identifying and Mitigating Collinearity

The collinearity problem is that the data do not contain enough "information" about the individual effects of explanatory variables to permit us to estimate all the parameters of the statistical model precisely

- Consequently, one solution is to obtain more information and include it in the analysis.
- A second way of adding new information is to introduce nonsample information in the form of restrictions on the parameters

6.5 Prediction

Page 86

Consider the model:

Eq. 6.32

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + e$$

- The prediction problem is to predict the value of the dependent variable y_0 , which is given by:

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

– The best linear unbiased predictor is:

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

Eq. 6.33

The variance of the forecast error, $f = (y_0 - \hat{y}_0)$, is: $\operatorname{var}(f) = \operatorname{var}[(\beta_1 + \beta_2 x_{02} + \beta_3 x_{03} + e_0) - (b_1 + b_2 x_{02} + b_3 x_{03})]$ $= \operatorname{var}(e_0 - b_1 - b_2 x_{02} - b_3 x_{03})$ $= \operatorname{var}(e_0) + \operatorname{var}(b_1) + x_{02}^2 \operatorname{var}(b_2) + x_{03}^2 \operatorname{var}(b_3)$ $+ 2x_{02} \operatorname{cov}(b_1, b_2) + 2x_{03} \operatorname{cov}(b_1, b_3) + 2x_{02} x_{03} \operatorname{cov}(b_2, b_3)$ 6.5 Prediction

6.5.1 An Example

> For our example, suppose $PRICE_0 = 6$, $ADVERT_0 = 1.9$, and $ADVERT_0^2 = 3.61$:

 $\begin{aligned} SALES_{0} &= 109.719 - 7.640 PRICE_{0} + 12.1512 ADVERT_{0} - 2.768 ADVERT_{0}^{2} \\ &= 109.719 - 7.640 \times 6 + 12.1512 \times 1.9 - 2.768 \times 3.61 \\ &= 76.974 \end{aligned}$

- We forecast sales will be \$76,974

6.5 Prediction	Table 6.3 Covariance Matrix for Andy's Burger Barn Model				
6.5.1 An Example					
		b_1	b_2	b_3	b_4
	b_1	46.227019	-6.426113	-11.600960	2.939026
	b_2	-6.426113	1.093988	0.300406	-0.085619
	b_3	-11.600960	0.300406	12.646302	-3.288746
	b_4	2.939026	-0.085619	-3.288746	0.884774

6.5.1 An Example

The estimated variance of the forecast error is: $\tilde{var}(f) = \hat{\sigma}^2 + \tilde{var}(b_1) + x_{02}^2 \tilde{var}(b_2) + x_{03}^2 \tilde{var}(b_3) + x_{04}^2 \tilde{var}(b_4)$ $+2x_{02}\overline{\text{cov}}(b_1, b_2) + 2x_{03}\overline{\text{cov}}(b_1, b_3) + 2x_{04}\overline{\text{cov}}(b_1, b_4)$ $+2x_{02}x_{03}\overline{\text{cov}}(b_2, b_3)+2x_{02}x_{04}\overline{\text{cov}}(b_2, b_4)+2x_{03}x_{04}\overline{\text{cov}}(b_3, b_4)$ $= 21.57865 + 46.22702 + 6^{2} \times 1.093988 + 1.9^{2} \times 12.6463 + 3.61^{2} \times 0.884774$ $+2 \times 6 \times (-6.426113) + 2 \times 1.9 \times (-11.60096) + 2 \times 3.61 \times 2.939026$ $+2 \times 6 \times 1.9 \times 0.300407 + 2 \times 6 \times 3.61 \times (-0.085619)$ $+2 \times 1.9 \times 3.61 \times (-3.288746)$ = 22.4208

– The standard error of the forecast error is:

$$\operatorname{se}(f) = \sqrt{22.4208} = 4.7351$$



6.5.1 An Example

The 95% prediction interval is:

 $(76.974 - 1.9939 \times 4.7351, 76.974 + 1.9939 \times 4.7351) = (67.533, 86.415)$

We predict, with 95% confidence, that the settings for price and advertising expenditure will yield *SALES* between \$67,533 and \$86,415

6.5 Prediction

6.5.1 An Example

The point forecast and the point estimate are both the same:

$$SALES_0 = E(SALES_0) = 76.974$$

– But:

$$se(E(SALES_0)) = \sqrt{var(f) - \hat{\sigma}^2} = \sqrt{22.4208 - 21.5786} = 0.9177$$

– A 95% confidence interval for $E(SALES_0)$ is:

 $(76.974 - 1.9939 \times 0.9177, 76.974 + 1.9939 \times 0.9177) = (75.144, 78.804)$

- AIC
 auxiliary regression
- BIC
- collinearity
- F-test
- irrelevantvariables
- nonsample information
- omitted variables

- omitted variable bias
- overall significance
- prediction
- RESET
 - restricted least squares
- restricted model
- restricted SSE
- SC

- single and joint null hypothesis
- testing many parameters
- unrestricted model
 - unrestricted SSE

Key Words

Appendices

Page 96

6A Chi-Square and *F*-Tests: More Details

Eq. 6A.1

Eq. 6A.2

The *F*-statistic is defined as:
 F = $\frac{(SSE_R - SSE_U)/J}{SSE_U/(N-K)}$ We can also show that:

$$V_1 = \frac{(SSE_R - SSE_U)}{\sigma^2} : \chi^2_{(J)}$$

For sufficiently large sample:

Eq. 6A.3

$$\hat{V}_1 = \frac{(SSE_R - SSE_U)}{\hat{\sigma}^2} : \chi^2_{(J)}$$

6A Chi-Square and *F*-Tests: More Details

Eq. 6A.4

But we can also show that:

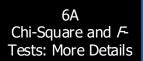
$$V_2 = \frac{(N-K)\hat{\sigma}^2}{\sigma^2} : \chi^2_{(N-K)}$$

From the book's appendix, we know that:

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

Eq. 6A.5 Therefore:

$$\frac{\frac{(SSE_R - SSE_U)}{\sigma^2}}{\frac{(N-K)\hat{\sigma}^2}{\sigma^2}} = \frac{[SSE_R - SSE_U]/J}{\hat{\sigma}^2} : F_{(J,N-K)}$$



A little reflection shows that:

$$F = \frac{\hat{V_1}}{J}$$

6A Chi-Square and *F*-Tests: More Details

When testing

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

in the equation

 $SALES = \beta_1 + \beta_2 PRICE + \beta_3 ADVERT + \beta_4 ADVERT^2 + e_i$

we get

$$F = 8.44$$
 p-value = .0005
 $\chi^2 = 16.88$ p-value = .0002

6A Chi-Square and *F*-Tests: More Details

Testing

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

we get

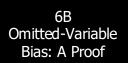
$$F = .936$$
 p-value = .3365
 $\chi^2 = .936$ p-value = .3333

Consider the model:

$$y_{i} = \beta_{1} + \beta_{2}x_{i2} + \beta_{3}x_{i3} + e_{i}$$

- Now suppose we incorrectly omit x_{i3} and estimate:

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



Notice the new disturbance term It's

$$v_i = \beta_3 x_{i3} + e_i$$

6B Omitted-Variable Bias: A Proof

The estimator for β_2 is:

Eq. 6B.1

$$b_{2}^{*} = \frac{\sum (x_{i2} - \overline{x}_{2})(y_{i} - \overline{y})}{\sum (x_{i2} - \overline{x}_{2})^{2}} = \beta_{2} + \sum w_{i}v_{i}$$

where

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

6B Omitted-Variable Bias: A Proof

Substituting for v_i yields:

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

where

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

6B Omitted-Variable Bias: A Proof

Hence:
$$E(b_{2}^{*}) = \beta_{2} + \beta_{3} \sum w_{i} x_{i3}$$

 $= \beta_{2} + \beta_{3} \frac{\sum (x_{i2} - \overline{x}_{2}) x_{i3}}{\sum (x_{i2} - \overline{x}_{2})^{2}}$
 $= \beta_{2} + \beta_{3} \frac{\sum (x_{i2} - \overline{x}_{2}) (x_{i3} - \overline{x}_{3})}{\sum (x_{i2} - \overline{x}_{2})^{2}}$
 $= \beta_{2} + \beta_{3} \frac{\overline{\text{cov}}(x_{2}, x_{3})}{\overline{\text{var}}(x_{2})} \neq \beta_{2}$