
Chapter 2

Properties of the regression coefficients and hypothesis testing

2.1 Overview

Chapter 1 introduced least squares regression analysis, a mathematical technique for fitting a relationship given suitable data on the variables involved. It is a fundamental chapter because much of the rest of the text is devoted to extending the least squares approach to handle more complex models, for example models with multiple explanatory variables, nonlinear models, and models with qualitative explanatory variables.

However, the mechanics of fitting regression equations are only part of the story. We are equally concerned with assessing the performance of our regression techniques and with developing an understanding of why they work better in some circumstances than in others. Chapter 2 is the starting point for this objective and is thus equally fundamental. In particular, it shows how two of the three main criteria for assessing the performance of estimators, unbiasedness and efficiency, are applied in the context of a regression model. The third criterion, consistency, will be considered in Chapter 8.

2.2 Learning outcomes

After working through the corresponding chapter in the text, studying the corresponding slideshows, and doing the starred exercises in the text and the additional exercises in this subject guide, you should be able to explain what is meant by:

- cross-sectional, time series, and panel data
- unbiasedness of OLS regression estimators
- variance and standard errors of regression coefficients and how they are determined
- Gauss–Markov theorem and efficiency of OLS regression estimators
- two-sided t tests of hypotheses relating to regression coefficients and one-sided t tests of hypotheses relating to regression coefficients
- F tests of goodness of fit of a regression equation

in the context of a regression model. The chapter is a long one and you should take your time over it because it is essential that you develop a perfect understanding of every detail.

2. Properties of the regression coefficients and hypothesis testing

2.3 Further material

Derivation of the expression for the variance of the naïve estimator in Section 2.3.

The variance of the naïve estimator in Section 2.3 and Exercise 2.9 is not of any great interest in itself, but its derivation provides an example of how one obtains expressions for variances of estimators in general.

In Section 2.3 we considered the naïve estimator of the slope coefficient derived by joining the first and last observations in a sample and calculating the slope of that line:

$$\hat{\beta}_2 = \frac{Y_n - Y_1}{X_n - X_1}.$$

It was demonstrated that the estimator could be decomposed as:

$$\hat{\beta}_2 = \beta_2 + \frac{u_n - u_1}{X_n - X_1}$$

and hence that $E(\hat{\beta}_2) = \beta_2$.

The population variance of a random variable X is defined to be $E([X - \mu_X]^2)$ where $\mu_X = E(X)$. Hence the population variance of $\hat{\beta}_2$ is given by:

$$\sigma_{\hat{\beta}_2}^2 = E([\hat{\beta}_2 - \beta_2]^2) = E\left(\left[\left\{\beta_2 + \frac{u_n - u_1}{X_n - X_1}\right\} - \beta_2\right]^2\right) = E\left(\left[\frac{u_n - u_1}{X_n - X_1}\right]^2\right).$$

On the assumption that X is nonstochastic, this can be written as:

$$\sigma_{\hat{\beta}_2}^2 = \left[\frac{1}{X_n - X_1}\right]^2 E([u_n - u_1]^2).$$

Expanding the quadratic, we have:

$$\begin{aligned} \sigma_{\hat{\beta}_2}^2 &= \left[\frac{1}{X_n - X_1}\right]^2 E(u_n^2 + u_1^2 - 2u_n u_1) \\ &= \left[\frac{1}{X_n - X_1}\right]^2 [E(u_n^2) + E(u_1^2) - 2E(u_n u_1)]. \end{aligned}$$

Each value of the disturbance term is drawn randomly from a distribution with mean 0 and population variance σ_u^2 , so $E(u_n^2)$ and $E(u_1^2)$ are both equal to σ_u^2 . u_n and u_1 are drawn independently from the distribution, so $E(u_n u_1) = E(u_n)E(u_1) = 0$. Hence:

$$\sigma_{\hat{\beta}_2}^2 = \frac{2\sigma_u^2}{(X_n - X_1)^2} = \frac{\sigma_u^2}{\frac{1}{2}(X_n - X_1)^2}.$$

Define $A = \frac{1}{2}(X_1 + X_n)$, the average of X_1 and X_n , and $D = X_n - A = A - X_1$. Then:

$$\begin{aligned}
 \frac{1}{2}(X_n - X_1)^2 &= \frac{1}{2}(X_n - A + A - X_1)^2 \\
 &= \frac{1}{2}[(X_n - A)^2 + (A - X_1)^2 + 2(X_n - A)(A - X_1)] \\
 &= \frac{1}{2}[D^2 + D^2 + 2(D)(D)] = 2D^2 \\
 &= (X_n - A)^2 + (A - X_1)^2 \\
 &= (X_n - A)^2 + (X_1 - A)^2 \\
 &= (X_n - \bar{X} + \bar{X} - A)^2 + (X_1 - \bar{X} + \bar{X} - A)^2 \\
 &= (X_n - \bar{X})^2 + (\bar{X} - A)^2 + 2(X_n - \bar{X})(\bar{X} - A) \\
 &\quad + (X_1 - \bar{X})^2 + (\bar{X} - A)^2 + 2(X_1 - \bar{X})(\bar{X} - A) \\
 &= (X_1 - \bar{X})^2 + (X_n - \bar{X})^2 + 2(\bar{X} - A)^2 + 2(X_1 + X_n - 2\bar{X})(\bar{X} - A) \\
 &= (X_1 - \bar{X})^2 + (X_n - \bar{X})^2 + 2(\bar{X} - A)^2 + 2(2A - 2\bar{X})(\bar{X} - A) \\
 &= (X_1 - \bar{X})^2 + (X_n - \bar{X})^2 - 2(\bar{X} - A)^2 \\
 &= (X_1 - \bar{X})^2 + (X_n - \bar{X})^2 - 2(A - \bar{X})^2 \\
 &= (X_1 - \bar{X})^2 + (X_n - \bar{X})^2 - \frac{1}{2}(X_1 + X_n - 2\bar{X})^2.
 \end{aligned}$$

Hence we obtain the expression in Exercise 2.9. There must be a shorter proof.

2.4 Additional exercises

A2.1 A variable Y depends on a nonstochastic variable X with the relationship:

$$Y = \beta_1 + \beta_2 X + u$$

where u is a disturbance term that satisfies the regression model assumptions. Given a sample of n observations, a researcher decides to estimate β_2 using the expression:

$$\hat{\beta}_2 = \frac{\sum X_i Y_i}{\sum X_i^2}.$$

(This is the OLS estimator of β_2 for the model $Y = \beta_2 X + u$.)

- Demonstrate that $\hat{\beta}_2$ is in general a biased estimator of β_2 .
- Discuss whether it is possible to determine the sign of the bias.
- Demonstrate that $\hat{\beta}_2$ is unbiased if $\beta_1 = 0$.
- Demonstrate that $\hat{\beta}_2$ is unbiased if $\bar{X} = 0$.

A2.2 A variable Y_i is generated as:

$$Y_i = \beta_1 + u_i$$

2. Properties of the regression coefficients and hypothesis testing

where β_1 is a fixed parameter and u_i is a disturbance term that is independently and identically distributed with expected value 0 and population variance σ_u^2 . The least squares estimator of β_1 is \bar{Y} , the sample mean of Y . However, a researcher believes that Y is a linear function of another variable X and uses ordinary least squares to fit the relationship:

$$\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X$$

calculating $\hat{\beta}_1$ as $\bar{Y} - \hat{\beta}_2 \bar{X}$, where \bar{X} is the sample mean of X . X may be assumed to be a nonstochastic variable. Determine whether the researcher's estimator $\hat{\beta}_1$ is biased or unbiased, and if biased, determine the direction of the bias.

A2.3 With the model described in Exercise A2.2, standard theory states that the population variance of the researcher's estimator of β_1 is:

$$\sigma_u^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right].$$

In general, this is larger than the population variance of \bar{Y} , which is σ_u^2/n . Explain the implications of the difference in the variances.

In the special case where $\bar{X} = 0$, the variances are the same. Give an intuitive explanation.

A2.4 A variable Y depends on a nonstochastic variable X with the relationship:

$$Y = \beta_1 + \beta_2 X + u$$

where u is a disturbance term that satisfies the regression model assumptions. Given a sample of n observations, a researcher decides to estimate β_2 using the expression:

$$\hat{\beta}_2 = \frac{\sum X_i Y_i}{\sum X_i^2}.$$

It can be shown that the population variance of this estimator is $\sigma_u^2 / \sum X_i^2$.

We saw in Exercise A2.1 that $\hat{\beta}_2$ is in general a biased estimator of β_2 . However, if either $\beta_1 = 0$ or $\bar{X} = 0$, the estimator is unbiased. What can be said in this case about the efficiency of the estimator in these two cases, comparing it with the estimator:

$$\frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}?$$

Returning to the general case where $\beta_1 \neq 0$ and $\bar{X} \neq 0$, suppose that there is very little variation in X in the sample. Is it possible that $\hat{\beta}_2$ might be a better estimator than the OLS estimator?

A2.5 Using the output for the regression in Exercise A1.1, reproduced below, perform appropriate statistical tests.

2.4. Additional exercises

```
. reg FDHO EXP if FDHO>0
```

Source	SS	df	MS			
Model	972602566	1	972602566	Number of obs =	6334	
Residual	1.7950e+09	6332	283474.003	F(1, 6332) =	3431.01	
				Prob > F =	0.0000	
				R-squared =	0.3514	
				Adj R-squared =	0.3513	
				Root MSE =	532.42	
Total	2.7676e+09	6333	437006.15			

FDHO	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
EXP	.0627099	.0010706	58.57	0.000	.0606112	.0648086
_cons	369.4418	10.65718	34.67	0.000	348.5501	390.3334

A2.6 Using the output for your regression in Exercise A1.2, perform appropriate statistical tests.

A2.7 Using the output for the regression of weight in 2004 on height in Exercise 1.9, reproduced below, perform appropriate statistical tests.

```
. reg WEIGHT04 HEIGHT
```

Source	SS	df	MS			
Model	211309	1	211309	Number of obs =	500	
Residual	595389.95	498	1195.56215	F(1, 498) =	176.74	
				Prob > F =	0.0000	
				R-squared =	0.2619	
				Adj R-squared =	0.2605	
				Root MSE =	34.577	
Total	806698.95	499	1616.63116			

WEIGHT04	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
HEIGHT	5.073711	.381639	13.29	0.000	4.32389	5.823532
_cons	-177.1703	25.93501	-6.83	0.000	-228.1258	-126.2147

A2.8 Using the output for the regression of earnings on height in Exercise A1.4, reproduced below, perform appropriate statistical tests.

```
. reg EARNINGS HEIGHT
```

Source	SS	df	MS			
Model	1393.77592	1	1393.77592	Number of obs =	500	
Residual	75171.3726	498	150.946531	F(1, 498) =	9.23	
				Prob > F =	0.0025	
				R-squared =	0.0182	
				Adj R-squared =	0.0162	
				Root MSE =	12.286	
Total	76565.1485	499	153.437171			

EARNINGS	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
HEIGHT	.4087231	.1345068	3.04	0.003	.1444523	.6729938
_cons	-9.26923	9.125089	-1.02	0.310	-27.19765	8.659188

2. Properties of the regression coefficients and hypothesis testing

- A2.9 Explain whether it would be justifiable to use a one-sided test on the slope coefficient in the regression of the rate of growth of employment on the rate of growth of GDP in Exercise 2.20.
- A2.10 Explain whether it would be justifiable to use a one-sided test on the slope coefficient in the regression of weight on height in Exercise 1.9.
- A2.11 With the information given in Exercise A1.5, how would the change in the measurement of GDP affect:
- the standard error of the coefficient of GDP
 - the F statistic for the equation?
- A2.12 With the information given in Exercise A1.6, how would the change in the measurement of GDP affect:
- the standard error of the coefficient of GDP
 - the F statistic for the equation?
- A2.13 [This is a continuation of Exercise 1.16 in the text.] A sample of data consists of n observations on two variables, Y and X . The true model is:

$$Y_i = \beta_1 + \beta_2 X_i + u_i$$

where β_1 and β_2 are parameters and u is a disturbance term that satisfies the usual regression model assumptions. In view of the true model:

$$\bar{Y} = \beta_1 + \beta_2 \bar{X} + \bar{u}$$

where \bar{Y} , \bar{X} , and \bar{u} are the sample means of Y , X , and u . Subtracting the second equation from the first, one obtains:

$$Y_i^* = \beta_2 X_i^* + u_i^*$$

where $Y_i^* = Y_i - \bar{Y}$, $X_i^* = X_i - \bar{X}$ and $u_i^* = u_i - \bar{u}$. Note that, by construction, the sample means of Y^* , X^* , and u^* are all equal to zero.

One researcher fits:

$$\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X. \quad (1)$$

A second researcher fits:

$$\hat{Y}^* = \hat{\beta}_1^* + \hat{\beta}_2^* X^*. \quad (2)$$

[Note: The second researcher included an intercept in the specification.]

- Comparing regressions (1) and (2), demonstrate that $\hat{Y}_i^* = \hat{Y}_i - \bar{Y}$.
- Demonstrate that the residuals in (2) are identical to the residuals in (1).
- Demonstrate that the OLS estimator of the variance of the disturbance term in (2) is equal to that in (1).
- Explain how the standard error of the slope coefficient in (2) is related to that in (1).

- Explain how R^2 in (2) is related to R^2 in (1).
- Explain why, theoretically, the specification (2) of the second researcher is incorrect and he should have fitted:

$$\widehat{Y}^* = \widehat{\beta}_2^* X^* \quad (3)$$

not including a constant in his specification.

- If the second researcher had fitted (3) instead of (2), how would this have affected his estimator of β_2 ? Would dropping the unnecessary intercept lead to a gain in efficiency?

A2.14 For the model described in Exercise A1.7, show that $\widehat{Y}_i^* = (\widehat{Y}_i - \bar{Y})/\widehat{\sigma}_Y$, and thus that $\widehat{u}_i^* = \widehat{u}_i/\widehat{\sigma}_Y$, where \widehat{Y}_i^* and \widehat{u}_i^* are the fitted value of Y_i^* and the residual in the transformed model.

Hence show that:

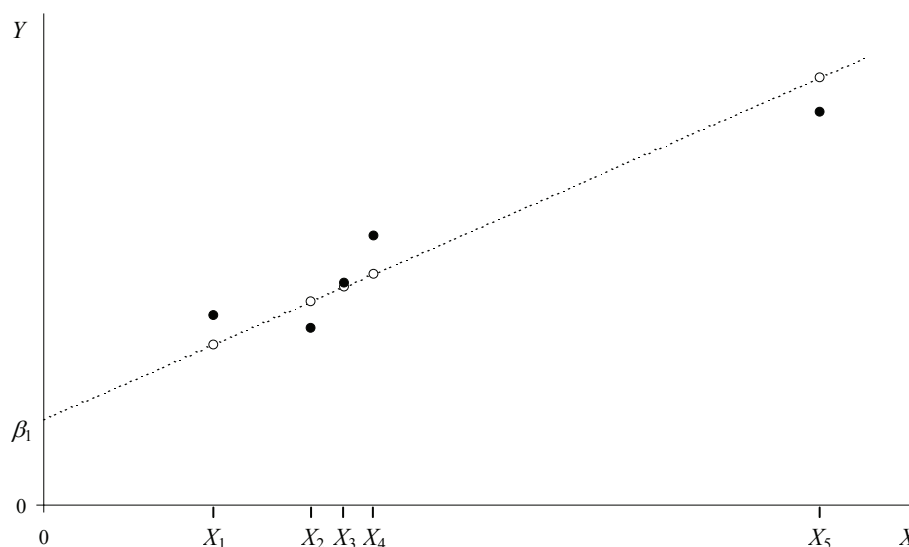
$$\text{s.e.}(\widehat{\beta}_2^*) = \frac{\widehat{\sigma}_X}{\widehat{\sigma}_Y} \times \text{s.e.}(\widehat{\beta}_2).$$

Hence find the relationship between the t statistic for $\widehat{\beta}_2^*$ and the t statistic for $\widehat{\beta}_2$ and the relationship between R^2 for the original specification and R^2 for the revised specification.

A2.15 A variable Y_i is generated as:

$$Y_i = \beta_1 + \beta_2 X_i + u_i \quad (1)$$

where β_1 and β_2 are fixed parameters and u_i is a disturbance term that satisfies the regression model assumptions. The values of X are fixed and are as shown in the figure. Four of them, X_1 to X_4 , are close together. The fifth, X_5 , is much larger. The corresponding values that Y would take, if there were no disturbance term, are given by the circles on the line. The presence of the disturbance term in the model causes the actual values of Y in a sample to be different. The solid black circles depict a typical sample of observations.



2. Properties of the regression coefficients and hypothesis testing

Discuss the advantages and disadvantages of dropping the observation corresponding to X_5 when regressing Y on X . If you keep the observation in the sample, will this cause the regression estimates to be biased?

2.5 Answers to the starred exercises in the textbook

2.1 Derive the decomposition of $\hat{\beta}_1$ shown in equation (2.29):

$$\hat{\beta}_1 = \beta_1 + \sum c_i u_i$$

where $c_i = \frac{1}{n} - a_i \bar{X}$ and a_i is defined in equation (2.23).

Answer:

$$\begin{aligned} \hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{X} &= (\beta_1 + \beta_2 \bar{X} + \bar{u}) - \bar{X} (\beta_2 + \sum a_i u_i) \\ &= \beta_1 + \frac{1}{n} \sum u_i - \bar{X} \sum a_i u_i \\ &= \beta_1 + \sum c_i u_i. \end{aligned}$$

2.5 An investigator correctly believes that the relationship between two variables X and Y is given by:

$$Y_i = \beta_1 + \beta_2 X_i + u_i.$$

Given a sample of observations on Y , X , and a third variable Z (which is not a determinant of Y), the investigator estimates β_2 as:

$$\frac{\sum (Z_i - \bar{Z})(Y_i - \bar{Y})}{\sum (Z_i - \bar{Z})(X_i - \bar{X})}.$$

Demonstrate that this estimator is unbiased.

Answer:

Noting that $Y_i - \bar{Y} = \beta_2 (X_i - \bar{X}) + u_i - \bar{u}$, we have:

$$\begin{aligned} \hat{\beta}_2 &= \frac{\sum (Z_i - \bar{Z})(Y_i - \bar{Y})}{\sum (Z_i - \bar{Z})(X_i - \bar{X})} \\ &= \frac{\sum (Z_i - \bar{Z}) \beta_2 (X_i - \bar{X}) + \sum (Z_i - \bar{Z})(u_i - \bar{u})}{\sum (Z_i - \bar{Z})(X_i - \bar{X})} \\ &= \beta_2 + \frac{\sum (Z_i - \bar{Z})(u_i - \bar{u})}{\sum (Z_i - \bar{Z})(X_i - \bar{X})}. \end{aligned}$$

2.5. Answers to the starred exercises in the textbook

Hence:

$$E(\hat{\beta}_2) = \beta_2 + \frac{\sum (Z_i - \bar{Z}) E(u_i - \bar{u})}{\sum (Z_i - \bar{Z})(X_i - \bar{X})} = \beta_2.$$

- 2.8 Using the decomposition of $\hat{\beta}_1$ obtained in Exercise 2.1, derive the expression for $\sigma_{\hat{\beta}_1}^2$ given in equation (2.42).

Answer:

$\hat{\beta}_1 = \beta_1 + \sum c_i u_i$, where $c_i = \frac{1}{n} - a_i \bar{X}$, and $E(\hat{\beta}_1) = \beta_1$. Hence:

$$\sigma_{\hat{\beta}_1}^2 = E \left[\left(\sum c_i u_i \right)^2 \right] = \sigma_u^2 \sum c_i^2 = \sigma_u^2 \left(n \frac{1}{n^2} - 2 \frac{\bar{X}}{n} \sum a_i + \bar{X}^2 \sum a_i^2 \right).$$

From Box 2.2, $\sum a_i = 0$ and:

$$\sum a_i^2 = \frac{1}{\sum (X_i - \bar{X})^2}.$$

Hence:

$$\sigma_{\hat{\beta}_1}^2 = \sigma_u^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right).$$

- 2.9 Given the decomposition in Exercise 2.2 of the OLS estimator of β_2 in the model $Y_i = \beta_2 X_i + u_i$, demonstrate that the variance of the slope coefficient is given by:

$$\sigma_{\hat{\beta}_2}^2 = \frac{\sigma_u^2}{\sum X_j^2}.$$

Answer:

$\hat{\beta}_2 = \beta_2 + \sum_{i=1}^n d_i u_i$, where $d_i = X_i / \sum_{j=1}^n X_j^2$, and $E(\hat{\beta}_2) = \beta_2$. Hence:

$$\begin{aligned} \sigma_{\hat{\beta}_2}^2 &= E \left[\left(\sum_{i=1}^n d_i u_i \right)^2 \right] = \sigma_u^2 \sum_{i=1}^n d_i^2 = \sigma_u^2 \sum_{i=1}^n \left(\frac{X_i^2}{\left(\sum_{j=1}^n X_j^2 \right)^2} \right) \\ &= \frac{\sigma_u^2}{\left(\sum_{j=1}^n X_j^2 \right)^2} \sum_{i=1}^n X_i^2 \\ &= \frac{\sigma_u^2}{\sum_{j=1}^n X_j^2}. \end{aligned}$$

2. Properties of the regression coefficients and hypothesis testing

2.12 It can be shown that the variance of the estimator of the slope coefficient in Exercise 2.5:

$$\frac{\sum (Z_i - \bar{Z})(Y_i - \bar{Y})}{\sum (Z_i - \bar{Z})(X_i - \bar{X})}$$

is given by:

$$\sigma_{\hat{\beta}_2}^2 = \frac{\sigma_u^2}{\sum (X_i - \bar{X})^2} \times \frac{1}{r_{XZ}^2}$$

where r_{XZ} is the correlation between X and Z . What are the implications for the efficiency of the estimator?

Answer:

If Z happens to be an exact linear function of X , the population variance will be the same as that of the OLS estimator. Otherwise $1/r_{XZ}^2$ will be greater than 1, the variance will be larger, and so the estimator will be less efficient.

2.15 Suppose that the true relationship between Y and X is $Y_i = \beta_1 + \beta_2 X_i + u_i$ and that the fitted model is $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$. In Exercise 1.13 it was shown that if $X_i^* = \mu_2 X_i$, and Y is regressed on X^* , the slope coefficient $\hat{\beta}_2^* = \hat{\beta}_2 / \mu_2$. How will the standard error of $\hat{\beta}_2^*$ be related to the standard error of $\hat{\beta}_2$?

Answer:

In Exercise 1.23 it was demonstrated that the fitted values of Y would be the same. This means that the residuals are the same, and hence $\hat{\sigma}_u^2$, the estimator of the variance of the disturbance term, is the same. The standard error of $\hat{\beta}_2^*$ is then given by:

$$\begin{aligned} \text{s.e.}(\hat{\beta}_2^*) &= \sqrt{\frac{\hat{\sigma}_u^2}{\sum (X_i^* - \bar{X}^*)^2}} \\ &= \sqrt{\frac{\hat{\sigma}_u^2}{\sum (\mu_2 X_i - \mu_2 \bar{X})^2}} \\ &= \sqrt{\frac{\hat{\sigma}_u^2}{\mu_2^2 \sum (X_i - \bar{X})^2}} \\ &= \frac{1}{\mu_2} \text{s.e.}(\hat{\beta}_2). \end{aligned}$$

2.17 A researcher with a sample of 50 individuals with similar education, but differing amounts of training, hypothesises that hourly earnings, *EARNINGS*, may be related to hours of training, *TRAINING*, according to the relationship:

$$EARNINGS = \beta_1 + \beta_2 TRAINING + u.$$

2.5. Answers to the starred exercises in the textbook

He is prepared to test the null hypothesis $H_0 : \beta_2 = 0$ against the alternative hypothesis $H_1 : \beta_2 \neq 0$ at the 5 per cent and 1 per cent levels. What should he report:

- (a) if $\hat{\beta}_2 = 0.30$, $\text{s.e.}(\hat{\beta}_2) = 0.12$?
- (b) if $\hat{\beta}_2 = 0.55$, $\text{s.e.}(\hat{\beta}_2) = 0.12$?
- (c) if $\hat{\beta}_2 = 0.10$, $\text{s.e.}(\hat{\beta}_2) = 0.12$?
- (d) if $\hat{\beta}_2 = -0.27$, $\text{s.e.}(\hat{\beta}_2) = 0.12$?

Answer:

There are 48 degrees of freedom, and hence the critical values of t at the 5 per cent, 1 per cent, and 0.1 per cent levels are 2.01, 2.68, and 3.51, respectively.

- (a) The t statistic is 2.50. Reject H_0 at the 5 per cent level but not at the 1 per cent level.
- (b) $t = 4.58$. Reject at the 0.1 per cent level.
- (c) $t = 0.83$. Fail to reject at the 5 per cent level.
- (d) $t = -2.25$. Reject H_0 at the 5 per cent level but not at the 1 per cent level.

2.22 Explain whether it would have been possible to perform one-sided tests instead of two-sided tests in Exercise 2.17. If you think that one-sided tests are justified, perform them and state whether the use of a one-sided test makes any difference.

Answer:

First, there should be a discussion of whether the parameter β_2 in:

$$EARNINGS = \beta_1 + \beta_2 TRAINING + u$$

can be assumed not to be negative. The objective of training is to impart skills. It would be illogical for an individual with greater skills to be paid less on that account, and so we can argue that we can rule out $\beta_2 < 0$. We can then perform a one-sided test. With 48 degrees of freedom, the critical values of t at the 5 per cent, 1 per cent, and 0.1 per cent levels are 1.68, 2.40, and 3.26, respectively.

- (a) The t statistic is 2.50. We can now reject H_0 at the 1 per cent level (but not at the 0.1 per cent level).
- (b) $t = 4.58$. Not affected by the change. Reject at the 0.1 per cent level.
- (c) $t = 0.83$. Not affected by the change. Fail to reject at the 5 per cent level.
- (d) $t = -2.25$. Reject H_0 at the 5 per cent level but not at the 1 per cent level.

Here there is a problem because the coefficient has an unexpected sign and is large enough to reject H_0 at the 5 per cent level with a two-sided test.

In principle we should ignore this and fail to reject H_0 . Admittedly, the likelihood of such a large negative t statistic occurring under H_0 is very small, but it would be smaller still under the alternative hypothesis $H_1 : \beta_2 > 0$.

However, we should consider two further possibilities. One is that the justification for a one-sided test is incorrect. For example, some jobs pay relatively low wages because they offer training that is valued by the employee.

2. Properties of the regression coefficients and hypothesis testing

Apprenticeships are the classic example. Alternatively, workers in some low-paid occupations may, for technical reasons, receive a relatively large amount of training. In either case, the correlation between training and earnings might be negative instead of positive.

Another possible reason for a coefficient having an unexpected sign is that the model is misspecified in some way. For example, the coefficient might be distorted by omitted variable bias, to be discussed in Chapter 6.

- 2.27 Suppose that the true relationship between Y and X is $Y_i = \beta_1 + \beta_2 X_i + u_i$ and that the fitted model is $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$. In Exercise 1.13 it was shown that if $X_i^* = \mu_2 X_i$, and Y is regressed on X^* , the slope coefficient $\hat{\beta}_2^* = \hat{\beta}_2 / \mu_2$. How will the t statistic for $\hat{\beta}_2^*$ be related to the t statistic for $\hat{\beta}_2$? (See also Exercise 2.15.)

Answer:

In Exercise 2.15 it was shown that $\text{s.e.}(\hat{\beta}_2^*) = \text{s.e.}(\hat{\beta}_2) / \mu_2$. Hence the t statistic is unaffected by the transformation.

Alternatively, since we saw in Exercise 1.23 that R^2 must be the same, it follows that the F statistic for the equation must be the same. For a simple regression the F statistic is the square of the t statistic on the slope coefficient, so the t statistic must be the same.

- 2.30 Calculate the 95 per cent confidence interval for β_2 in the price inflation/wage inflation example:

$$\hat{p} = \begin{matrix} -1.21 & + & 0.82w. \\ (0.05) & & (0.10) \end{matrix}$$

What can you conclude from this calculation?

Answer:

With n equal to 20, there are 18 degrees of freedom and the critical value of t at the 5 per cent level is 2.10. The 95 per cent confidence interval is therefore:

$$0.82 - 0.10 \times 2.10 \leq \beta_2 \leq 0.82 + 0.10 \times 2.10$$

that is:

$$0.61 \leq \beta_2 \leq 1.03.$$

We see that we cannot (quite) reject the null hypothesis $H_0 : \beta_2 = 1$.

- 2.36 Suppose that the true relationship between Y and X is $Y_i = \beta_1 + \beta_2 X_i + u_i$ and that the fitted model is $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$. Suppose that $X_i^* = \mu_2 X_i$, and Y is regressed on X^* . How will the F statistic for this regression be related to the F statistic for the original regression? (See also Exercises 1.23, 2.15, and 2.27.)

Answer:

We saw in Exercise 1.23 that R^2 would be the same, and it follows that F must also be the same.

2.6 Answers to the additional exercises

Note: Each of the exercises below relates to a simple regression. Accordingly, the F test is equivalent to a two-sided t test on the slope coefficient and there is no point in performing both tests. The F statistic is equal to the square of the t statistic and, for any significance level, the critical value of F is equal to the critical value of t . Obviously a one-sided t test, when justified, is preferable to either in that it has greater power for any given significance level.

A2.1 We have:

$$\hat{\beta}_2 = \frac{\sum X_i Y_i}{\sum X_i^2} = \frac{\sum X_i (\beta_1 + \beta_2 X_i + u_i)}{\sum X_i^2} = \frac{\beta_1 \sum X_i}{\sum X_i^2} + \beta_2 + \frac{\sum X_i u_i}{\sum X_i^2}.$$

Hence:

$$E(\hat{\beta}_2) = \frac{\beta_1 \sum X_i}{\sum X_i^2} + \beta_2 + E\left(\frac{\sum X_i u_i}{\sum X_i^2}\right) = \frac{\beta_1 \sum X_i}{\sum X_i^2} + \beta_2 + \frac{\sum X_i E(u_i)}{\sum X_i^2}$$

assuming that X is nonstochastic. Since $E(u_i) = 0$, then:

$$E(\hat{\beta}_2) = \frac{\beta_1 \sum X_i}{\sum X_i^2} + \beta_2.$$

Thus $\hat{\beta}_2$ will in general be a biased estimator. The sign of the bias depends on the signs of β_1 and $\sum X_i$. In general, we have no information about either of these. However, if either $\beta_1 = 0$ or $\bar{X} = 0$ (and so $\sum X_i = 0$), the bias term disappears and $\hat{\beta}_2$ is unbiased after all.

A2.2 First we need to show that $E(\hat{\beta}_2) = 0$.

$$\hat{\beta}_2 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2} = \frac{\sum (X_i - \bar{X})(\beta_1 + u_i - \beta_1 - \bar{u})}{\sum (X_i - \bar{X})^2} = \frac{\sum (X_i - \bar{X})(u_i - \bar{u})}{\sum (X_i - \bar{X})^2}.$$

Hence, given that we are told that X is nonstochastic:

$$\begin{aligned} E(\hat{\beta}_2) &= E\left(\frac{\sum (X_i - \bar{X})(u_i - \bar{u})}{\sum (X_i - \bar{X})^2}\right) \\ &= \frac{1}{\sum (X_i - \bar{X})^2} E\left(\sum (X_i - \bar{X})(u_i - \bar{u})\right) \\ &= \frac{1}{\sum (X_i - \bar{X})^2} \sum (X_i - \bar{X}) E(u_i - \bar{u}) \\ &= 0 \end{aligned}$$

since $E(u) = 0$. Thus:

$$E(\hat{\beta}_1) = E(\bar{Y} - \hat{\beta}_2 \bar{X}) = \beta_1 - \bar{X} E(\hat{\beta}_2) = \beta_1$$

and the estimator is unbiased.

2. Properties of the regression coefficients and hypothesis testing

A2.3 In general, the researcher's estimator will have a larger variance than \bar{Y} and therefore will be inefficient. However, if $\bar{X} = 0$, the variances are the same. This is because the estimators are then identical. $\bar{Y} - \hat{\beta}_2 \bar{X}$ reduces to \bar{Y} .

A2.4 The variance of the estimator is $\sigma_u^2 / \sum X_i^2$ whereas that of the estimator:

$$\frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

is:

$$\frac{\sigma_u^2}{\sum (X_i - \bar{X})^2} = \frac{\sigma_u^2}{\sum X_i^2 - n\bar{X}^2}.$$

Thus, provided $\bar{X} \neq 0$, $\sigma_u^2 / \sum X_i^2$ is more efficient than:

$$\frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

if $\beta_1 = 0$ because it is unbiased and has a smaller variance. It is the OLS estimator in this case.

If $\bar{X} = 0$, the estimators are equally efficient because the population variance expressions are identical. The reason for this is that the estimators are now identical:

$$\frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2} = \frac{\sum X_i(Y_i - \bar{Y})}{\sum X_i^2} = \frac{\sum X_i Y_i}{\sum X_i^2} - \frac{\bar{Y} \sum X_i}{\sum X_i^2} = \frac{\sum X_i Y_i}{\sum X_i^2}$$

since $\sum X_i = n\bar{X} = 0$.

Returning to the general case, if there is little variation in X in the sample, $\sum (X_i - \bar{X})^2$ may be small and hence the population variance of $\sum (X_i - \bar{X})(Y_i - \bar{Y}) / \sum (X_i - \bar{X})^2$ may be large. Thus using a criterion such as mean square error, $\hat{\beta}_2$ may be preferable if the bias is small.

A2.5 The t statistic for the coefficient of EXP is 58.57, very highly significant. There is little point performing a t test on the intercept, given that it has no plausible meaning. The F statistic is 3431.0, very highly significant. Since this is a simple regression model, the two tests are equivalent.

A2.6 The slope coefficient for every category is significantly different from zero at a very high significance level. (The F test is equivalent to the t test on the slope coefficient.)

2.6. Answers to the additional exercises

<i>EXP</i>						
	<i>n</i>	$\hat{\beta}_2$	s.e. ($\hat{\beta}_2$)	<i>t</i>	R^2	<i>F</i>
<i>ADM</i>	2,815	0.0235	0.0008	28.86	0.228	832.8
<i>CLOT</i>	4,500	0.0316	0.0010	30.99	0.176	960.6
<i>DOM</i>	1,661	0.0409	0.0026	16.02	0.134	256.6
<i>EDUC</i>	561	0.1202	0.0090	13.30	0.241	177.0
<i>ELEC</i>	5,828	0.0131	0.0004	35.70	0.180	1274.8
<i>FDAW</i>	5,102	0.0527	0.0010	52.86	0.354	2794.7
<i>FDHO</i>	6,334	0.0627	0.0011	58.57	0.351	3431.0
<i>FOOT</i>	1,827	0.0058	0.0005	12.78	0.082	163.4
<i>FURN</i>	487	0.0522	0.0070	7.44	0.102	55.3
<i>GASO</i>	5,710	0.0373	0.0008	46.89	0.278	2198.5
<i>HEAL</i>	4,802	0.0574	0.0018	31.83	0.174	1013.4
<i>HOUS</i>	6,223	0.1976	0.0027	74.16	0.469	5499.9
<i>LIFE</i>	1,253	0.0193	0.0016	11.86	0.101	140.7
<i>LOCT</i>	692	0.0068	0.0010	6.59	0.059	43.5
<i>MAPP</i>	399	0.0329	0.0049	6.72	0.102	45.1
<i>PERS</i>	3,817	0.0069	0.0002	32.15	0.213	1033.4
<i>READ</i>	2,287	0.0048	0.0003	16.28	0.104	265.1
<i>SAPP</i>	1,037	0.0045	0.0007	6.03	0.034	36.4
<i>TELE</i>	5,788	0.0160	0.0003	46.04	0.268	2119.7
<i>TEXT</i>	992	0.0040	0.0006	7.32	0.051	53.5
<i>TOB</i>	1,155	0.0165	0.0016	10.56	0.088	111.6
<i>TOYS</i>	2,504	0.0145	0.0010	14.34	0.076	205.7
<i>TRIP</i>	516	0.0466	0.0043	10.84	0.186	117.5

A2.7 The *t* statistic, 13.29, is very highly significant. (The *F* test is equivalent.)

A2.8 The *t* statistic for height, 3.04, suggests that the effect of height on earnings is highly significant, despite the very low R^2 . In principle the estimate of an extra 41 cents of hourly earnings for every extra inch of height could have been a purely random result of the kind that one obtains with nonsense models. However, the fact that it is apparently highly significant causes us to look for other explanations, the most likely one being that suggested in the answer to Exercise A1.4. Of course, we would not attempt to test the negative constant.

A2.9 One could justify a one-sided test on the slope coefficient in the regression of the rate of growth of employment on the rate of growth of GDP on the grounds that an increase in the rate of growth of GDP is unlikely to cause a decrease in the rate of growth of employment.

A2.10 One could justify a one-sided test on the slope coefficient in the regression of weight on height in Exercise 1.9 on the grounds that an increase in height is unlikely to cause a decrease in weight.

2. Properties of the regression coefficients and hypothesis testing

A2.11 *The standard error of the coefficient of GDP.* This is given by:

$$\frac{\sqrt{\widehat{\sigma}_u^{*2}}}{\sqrt{\sum (G_i^* - \bar{G}^*)^2}}$$

where $\widehat{\sigma}_u^{*2}$, the estimator of the variance of the disturbance term, is $\sum \widehat{u}_i^{*2}/(n-2)$. Since RSS is unchanged, $\widehat{\sigma}_u^{*2} = \widehat{\sigma}_u^2$.

We saw in Exercise A1.6 that $G_i^* - \bar{G}^* = G_i - \bar{G}$ for all i . Hence the new standard error is given by:

$$\frac{\sqrt{\widehat{\sigma}_u^2}}{\sqrt{\sum (G_i - \bar{G})^2}}$$

and is unchanged.

$$F = \frac{ESS}{RSS/(n-2)}$$

where:

$$ESS = \text{explained sum of squares} = \sum (\widehat{Y}_i^* - \bar{Y}^*)^2.$$

Since $\widehat{u}_i^* = \widehat{u}_i$, $\widehat{Y}_i^* = \widehat{Y}_i$ and ESS is unchanged. We saw in Exercise A1.6 that RSS is unchanged. Hence F is unchanged.

A2.12 *The standard error of the coefficient of GDP.* This is given by:

$$\frac{\sqrt{\widehat{\sigma}_u^{*2}}}{\sqrt{\sum (G_i^* - \bar{G}^*)^2}}$$

where $\widehat{\sigma}_u^{*2}$, the estimator of the variance of the disturbance term, is $\sum \widehat{u}_i^{*2}/(n-2)$. We saw in Exercise 1.7 that $\widehat{u}_i^* = \widehat{u}_i$ and so RSS is unchanged. Hence $\widehat{\sigma}_u^{*2} = \widehat{\sigma}_u^2$. Thus the new standard error is given by:

$$\frac{\sqrt{\widehat{\sigma}_u^2}}{\sqrt{\sum (2G_i - 2\bar{G})^2}} = \frac{1}{2} \frac{\sqrt{\widehat{\sigma}_u^2}}{\sqrt{\sum (G_i - \bar{G})^2}} = 0.005.$$

$F = ESS/(RSS/(n-2))$ where:

$$ESS = \text{explained sum of squares} = \sum (\widehat{Y}_i^* - \bar{Y}^*)^2.$$

Since $\widehat{u}_i^* = \widehat{u}_i$, $\widehat{Y}_i^* = \widehat{Y}_i$ and ESS is unchanged. Hence F is unchanged.

A2.13 One way of demonstrating that $\widehat{Y}_i^* = \widehat{Y}_i - \bar{Y}$:

$$\begin{aligned} \widehat{Y}_i^* &= \widehat{\beta}_1^* + \widehat{\beta}_2^* X_i^* = \widehat{\beta}_2 (X_i - \bar{X}) \\ \widehat{Y}_i - \bar{Y} &= (\widehat{\beta}_1 + \widehat{\beta}_2 X_i) - \bar{Y} = (\bar{Y} - \widehat{\beta}_2 \bar{X}) + \widehat{\beta}_2 X_i - \bar{Y} = \widehat{\beta}_2 (X_i - \bar{X}). \end{aligned}$$

Demonstration that the residuals are the same:

$$\hat{u}_i^* = Y_i^* - \hat{Y}_i^* = (Y_i - \bar{Y}) - (\hat{Y}_i - \bar{Y}) = \hat{u}_i.$$

Demonstration that the OLS estimator of the variance of the disturbance term in (2) is equal to that in (1):

$$\hat{\sigma}_u^{*2} = \frac{\sum \hat{u}_i^{*2}}{n-2} = \frac{\sum \hat{u}_i^2}{n-2} = \hat{\sigma}_u^2.$$

The standard error of the slope coefficient in (2) is equal to that in (1).

$$\hat{\sigma}_{\hat{\beta}_2^*}^2 = \frac{\hat{\sigma}_u^{*2}}{\sum (X_i^* - \bar{X})^2} = \frac{\hat{\sigma}_u^2}{\sum X_i^{*2}} = \frac{\hat{\sigma}_u^2}{\sum (X_i - \bar{X})^2} = \hat{\sigma}_{\hat{\beta}_2}^2.$$

Hence the standard errors are the same.

Demonstration that R^2 in (2) is equal to R^2 in (1):

$$R^{2*} = \frac{\sum (\hat{Y}_i^* - \bar{Y}^*)^2}{\sum (Y_i^* - \bar{Y}^*)^2}.$$

$\hat{Y}_i^* = \hat{Y}_i - \bar{Y}$ and $\bar{Y}^* = \bar{Y}$. Hence $\hat{Y}^* = 0$. $\bar{Y}^* = \bar{Y} - \bar{Y} = 0$. Hence:

$$R^{2*} = \frac{\sum (\hat{Y}_i^*)^2}{\sum (Y_i^*)^2} = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = R^2.$$

The reason that specification (2) of the second researcher is incorrect is that the model does not include an intercept.

If the second researcher had fitted (3) instead of (2), this would not in fact have affected his estimator of β_2 . Using (3), the researcher should have estimated β_2 as:

$$\hat{\beta}_2^* = \frac{\sum X_i^* Y_i^*}{\sum X_i^{*2}}.$$

However, Exercise 1.16 demonstrates that, effectively, he has done exactly this. Hence the estimator will be the same. It follows that dropping the unnecessary intercept would not have led to a gain in efficiency.

A2.14 We have:

$$\hat{Y}_i^* = \hat{\beta}_2^* X_i^* = \frac{\hat{\sigma}_X}{\hat{\sigma}_Y} \hat{\beta}_2 \left(\frac{X_i - \bar{X}}{\hat{\sigma}_X} \right) = \frac{1}{\hat{\sigma}_Y} \hat{\beta}_2 (X_i - \bar{X})$$

and:

$$\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i = (\bar{Y} - \hat{\beta}_2 \bar{X}) + \hat{\beta}_2 X_i = \bar{Y} + \hat{\beta}_2 (X_i - \bar{X}).$$

Hence:

$$\hat{Y}_i^* = \frac{1}{\hat{\sigma}_Y} (\hat{Y}_i - \bar{Y}).$$

2. Properties of the regression coefficients and hypothesis testing

Also:

$$\hat{u}_i^* = Y_i^* - \hat{Y}_i^* = \frac{1}{\hat{\sigma}_Y}(Y_i - \bar{Y}) - \frac{1}{\hat{\sigma}_Y}(\hat{Y}_i - \bar{Y}) = \frac{1}{\hat{\sigma}_Y}(Y_i - \hat{Y}_i) = \frac{1}{\hat{\sigma}_Y}\hat{u}_i$$

and:

$$\text{s.e.}(\hat{\beta}_2^*) = \sqrt{\frac{\frac{1}{n-2} \sum \hat{u}_i^{*2}}{\sum (X_i^* - \bar{X}^*)^2}} = \sqrt{\frac{\left(\frac{1}{\hat{\sigma}_Y}\right)^2 \frac{1}{n-2} \sum \hat{u}_i^2}{\sum \left(\frac{X_i - \bar{X}}{\hat{\sigma}_X}\right)^2}} = \frac{\hat{\sigma}_X}{\hat{\sigma}_Y} \times \text{s.e.}(\hat{\beta}_2).$$

Given the expressions for $\hat{\beta}_2^*$ and $\text{s.e.}(\hat{\beta}_2^*)$, the t statistic for $\hat{\beta}_2^*$ is the same as that for $\hat{\beta}_2$. Hence the F statistic will be the same and R^2 will be the same.

A2.15 The inclusion of the fifth observation does not cause the model to be misspecified or the regression model assumptions to be violated, so retaining it in the sample will not give rise to biased estimates. There would be no advantages in dropping it and there would be one major disadvantage. $\sum (X_i - \bar{X})^2$ would be greatly reduced and hence the variances of the coefficients would be increased, adversely affecting the precision of the estimates.

This said, in practice one would wish to check whether it is sensible to assume that the model relating Y to X for the other observations really does apply to the observation corresponding to X_5 as well. This question can be answered only by being familiar with the context and having some intuitive understanding of the relationship between Y and X .