

# Regression with a Single Regressor: Hypothesis Tests and Confidence Intervals

This chapter continues the treatment of linear regression with a single regressor. Chapter 4 explained how the OLS estimator  $\hat{\beta}_1$  of the slope coefficient  $\beta_1$  differs from one sample to the next—that is, how  $\hat{\beta}_1$  has a sampling distribution. In this chapter, we show how knowledge of this sampling distribution can be used to make statements about  $\beta_1$  that accurately summarize the sampling uncertainty. The starting point is the standard error of the OLS estimator, which measures the spread of the sampling distribution of  $\hat{\beta}_1$ . Section 5.1 provides an expression for this standard error (and for the standard error of the OLS estimator of the intercept), then shows how to use  $\hat{\beta}_1$  and its standard error to test hypotheses. Section 5.2 explains how to construct confidence intervals for  $\beta_1$ . Section 5.3 takes up the special case of a binary regressor.

Sections 5.1 through 5.3 assume that the three least squares assumptions of Chapter 4 hold. If, in addition, some stronger conditions hold, then some stronger results can be derived regarding the distribution of the OLS estimator. One of these stronger conditions is that the errors are homoskedastic, a concept introduced in Section 5.4. Section 5.5 presents the Gauss–Markov theorem, which states that, under certain conditions, OLS is efficient (has the smallest variance) among a certain class of estimators. Section 5.6 discusses the distribution of the OLS estimator when the population distribution of the regression errors is normal.

## 5.1 Testing Hypotheses About One of the Regression Coefficients

Your client, the superintendent, calls you with a problem. She has an angry taxpayer in her office who asserts that cutting class size will not help boost test scores, so reducing them further is a waste of money. Class size, the taxpayer claims, has no effect on test scores.

The taxpayer's claim can be rephrased in the language of regression analysis. Because the effect on test scores of a unit change in class size is  $\beta_{ClassSize}$ , the taxpayer is asserting that the population regression line is flat—that is, the slope  $\beta_{ClassSize}$  of the population regression line is zero. Is there, the superintendent asks,

## General Form of the $t$ -Statistic

**KEY CONCEPT**

5.1

In general, the  $t$ -statistic has the form

$$t = \frac{\text{estimator} - \text{hypothesized value}}{\text{standard error of the estimator}}. \quad (5.1)$$

evidence in your sample of 420 observations on California school districts that this slope is nonzero? Can you reject the taxpayer's hypothesis that  $\beta_{ClassSize} = 0$ , or should you accept it, at least tentatively pending further new evidence?

This section discusses tests of hypotheses about the slope  $\beta_1$  or intercept  $\beta_0$  of the population regression line. We start by discussing two-sided tests of the slope  $\beta_1$  in detail, then turn to one-sided tests and to tests of hypotheses regarding the intercept  $\beta_0$ .

### Two-Sided Hypotheses Concerning $\beta_1$

The general approach to testing hypotheses about the coefficient  $\beta_1$  is the same as to testing hypotheses about the population mean, so we begin with a brief review.

**Testing hypotheses about the population mean.** Recall from Section 3.2 that the null hypothesis that the mean of  $Y$  is a specific value  $\mu_{Y,0}$  can be written as  $H_0: E(Y) = \mu_{Y,0}$ , and the two-sided alternative is  $H_1: E(Y) \neq \mu_{Y,0}$ .

The test of the null hypothesis  $H_0$  against the two-sided alternative proceeds as in the three steps summarized in Key Concept 3.6. The first is to compute the standard error of  $\bar{Y}$ ,  $SE(\bar{Y})$ , which is an estimator of the standard deviation of the sampling distribution of  $\bar{Y}$ . The second step is to compute the  $t$ -statistic, which has the general form given in Key Concept 5.1; applied here, the  $t$ -statistic is  $t = (\bar{Y} - \mu_{Y,0})/SE(\bar{Y})$ .

The third step is to compute the  $p$ -value, which is the smallest significance level at which the null hypothesis could be rejected, based on the test statistic actually observed; equivalently, the  $p$ -value is the probability of obtaining a statistic, by random sampling variation, at least as different from the null hypothesis value as is the statistic actually observed, assuming that the null hypothesis is correct (Key Concept 3.5). Because the  $t$ -statistic has a standard normal distribution in large samples under the null hypothesis, the  $p$ -value for a two-sided hypothesis test is  $2\Phi(-|t^{act}|)$ , where  $t^{act}$  is the value of the  $t$ -statistic actually computed and  $\Phi$  is the cumulative standard normal distribution tabulated in Appendix Table 1. Alternatively,

the third step can be replaced by simply comparing the  $t$ -statistic to the critical value appropriate for the test with the desired significance level. For example, a two-sided test with a 5% significance level would reject the null hypothesis if  $|t^{act}| > 1.96$ . In this case, the population mean is said to be statistically significantly different from the hypothesized value at the 5% significance level.

**Testing hypotheses about the slope  $\beta_1$ .** At a theoretical level, the critical feature justifying the foregoing testing procedure for the population mean is that, in large samples, the sampling distribution of  $\bar{Y}$  is approximately normal. Because  $\hat{\beta}_1$  also has a normal sampling distribution in large samples, hypotheses about the true value of the slope  $\beta_1$  can be tested using the same general approach.

The null and alternative hypotheses need to be stated precisely before they can be tested. The angry taxpayer's hypothesis is that  $\beta_{ClassSize} = 0$ . More generally, under the null hypothesis the true population slope  $\beta_1$  takes on some specific value,  $\beta_{1,0}$ . Under the two-sided alternative,  $\beta_1$  does not equal  $\beta_{1,0}$ . That is, the **null hypothesis** and the **two-sided alternative hypothesis** are

$$H_0: \beta_1 = \beta_{1,0} \text{ vs. } H_1: \beta_1 \neq \beta_{1,0} \quad (\text{two-sided alternative}). \quad (5.2)$$

To test the null hypothesis  $H_0$ , we follow the same three steps as for the population mean.

The first step is to compute the **standard error of  $\hat{\beta}_1$** ,  $SE(\hat{\beta}_1)$ . The standard error of  $\hat{\beta}_1$  is an estimator of  $\sigma_{\hat{\beta}_1}$ , the standard deviation of the sampling distribution of  $\hat{\beta}_1$ . Specifically,

$$SE(\hat{\beta}_1) = \sqrt{\hat{\sigma}_{\hat{\beta}_1}^2}, \quad (5.3)$$

where

$$\hat{\sigma}_{\hat{\beta}_1}^2 = \frac{1}{n} \times \frac{\frac{1}{n-2} \sum_{i=1}^n (X_i - \bar{X})^2 \hat{u}_i^2}{\left[ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right]^2}. \quad (5.4)$$

The estimator of the variance in Equation (5.4) is discussed in Appendix 5.1. Although the formula for  $\hat{\sigma}_{\hat{\beta}_1}^2$  is complicated, in applications the standard error is computed by regression software so that it is easy to use in practice.

The second step is to compute the  **$t$ -statistic**,

$$t = \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)}. \quad (5.5)$$

Hypothesis tests are useful if you have a specific null hypothesis in mind (as did our angry taxpayer). Being able to accept or to reject this null hypothesis based on the statistical evidence provides a powerful tool for coping with the uncertainty inherent in using a sample to learn about the population. Yet, there are many times that no single hypothesis about a regression coefficient is dominant, and instead one would like to know a range of values of the coefficient that are consistent with the data. This calls for constructing a confidence interval.

## 5.2 Confidence Intervals for a Regression Coefficient

Because any statistical estimate of the slope  $\beta_1$  necessarily has sampling uncertainty, we cannot determine the true value of  $\beta_1$  exactly from a sample of data. It is possible, however, to use the OLS estimator and its standard error to construct a confidence interval for the slope  $\beta_1$  or for the intercept  $\beta_0$ .

**Confidence interval for  $\beta_1$ .** Recall that a 95% **confidence interval for  $\beta_1$**  has two equivalent definitions. First, it is the set of values that cannot be rejected using a two-sided hypothesis test with a 5% significance level. Second, it is an interval that has a 95% probability of containing the true value of  $\beta_1$ ; that is, in 95% of possible samples that might be drawn, the confidence interval will contain the true value of  $\beta_1$ . Because this interval contains the true value in 95% of all samples, it is said to have a **confidence level** of 95%.

The reason these two definitions are equivalent is as follows. A hypothesis test with a 5% significance level will, by definition, reject the true value of  $\beta_1$  in only 5% of all possible samples; that is, in 95% of all possible samples, the true value of  $\beta_1$  will *not* be rejected. Because the 95% confidence interval (as defined in the first definition) is the set of all values of  $\beta_1$  that are *not* rejected at the 5% significance level, it follows that the true value of  $\beta_1$  will be contained in the confidence interval in 95% of all possible samples.

As in the case of a confidence interval for the population mean (Section 3.3), in principle a 95% confidence interval can be computed by testing all possible values of  $\beta_1$  (that is, testing the null hypothesis  $\beta_1 = \beta_{1,0}$  for all values of  $\beta_{1,0}$ ) at the 5% significance level using the *t*-statistic. The 95% confidence interval is then the collection of all the values of  $\beta_1$  that are not rejected. But constructing the *t*-statistic for all values of  $\beta_1$  would take forever.

An easier way to construct the confidence interval is to note that the *t*-statistic will reject the hypothesized value  $\beta_{1,0}$  whenever  $\beta_{1,0}$  is outside the range

**KEY CONCEPT** **Confidence Interval for  $\beta_1$** **5.3**

A 95% two-sided confidence interval for  $\beta_1$  is an interval that contains the true value of  $\beta_1$  with a 95% probability; that is, it contains the true value of  $\beta_1$  in 95% of all possible randomly drawn samples. Equivalently, it is the set of values of  $\beta_1$  that cannot be rejected by a 5% two-sided hypothesis test. When the sample size is large, it is constructed as

$$\begin{aligned} \text{95\% confidence interval for } \beta_1 = \\ [\hat{\beta}_1 - 1.96SE(\hat{\beta}_1), \hat{\beta}_1 + 1.96SE(\hat{\beta}_1)]. \end{aligned} \quad (5.12)$$

$\hat{\beta}_1 \pm 1.96SE(\hat{\beta}_1)$ . That is, the 95% confidence interval for  $\beta_1$  is the interval  $[\hat{\beta}_1 - 1.96SE(\hat{\beta}_1), \hat{\beta}_1 + 1.96SE(\hat{\beta}_1)]$ . This argument parallels the argument used to develop a confidence interval for the population mean.

The construction of a confidence interval for  $\beta_1$  is summarized as Key Concept 5.3.

**Confidence interval for  $\beta_0$**  A 95% confidence interval for  $\beta_0$  is constructed as in Key Concept 5.3, with  $\hat{\beta}_0$  and  $SE(\hat{\beta}_0)$  replacing  $\hat{\beta}_1$  and  $SE(\hat{\beta}_1)$ .

**Application to test scores.** The OLS regression of the test score against the student-teacher ratio, reported in Equation (5.8), yielded  $\hat{\beta}_1 = -2.28$  and  $SE(\hat{\beta}_1) = 0.52$ . The 95% two-sided confidence interval for  $\beta_1$  is  $\{-2.28 \pm 1.96 \times 0.52\}$ , or  $-3.30 \leq \beta_1 \leq -1.26$ . The value  $\beta_1 = 0$  is not contained in this confidence interval, so (as we knew already from Section 5.1) the hypothesis  $\beta_1 = 0$  can be rejected at the 5% significance level.

**Confidence intervals for predicted effects of changing  $X$ .** The 95% confidence interval for  $\beta_1$  can be used to construct a 95% confidence interval for the predicted effect of a general change in  $X$ .

Consider changing  $X$  by a given amount,  $\Delta x$ . The predicted change in  $Y$  associated with this change in  $X$  is  $\beta_1 \Delta x$ . The population slope  $\beta_1$  is unknown, but because we can construct a confidence interval for  $\beta_1$ , we can construct a confidence interval for the predicted effect  $\beta_1 \Delta x$ . Because one end of a 95% confidence interval for  $\beta_1$  is  $\hat{\beta}_1 - 1.96SE(\hat{\beta}_1)$ , the predicted effect of the change  $\Delta x$  using this estimate of  $\beta_1$  is  $[\hat{\beta}_1 - 1.96SE(\hat{\beta}_1)] \times \Delta x$ . The other end of the confidence interval

is  $\hat{\beta}_1 + 1.96SE(\hat{\beta}_1)$ , and the predicted effect of the change using that estimate is  $[\hat{\beta}_1 + 1.96SE(\hat{\beta}_1)] \times \Delta x$ . Thus a 95% confidence interval for the effect of changing  $x$  by the amount  $\Delta x$  can be expressed as

$$\begin{aligned} & \text{95\% confidence interval for } \beta_1 \Delta x = \\ & [\hat{\beta}_1 \Delta x - 1.96SE(\hat{\beta}_1) \times \Delta x, \hat{\beta}_1 \Delta x + 1.96SE(\hat{\beta}_1) \times \Delta x]. \end{aligned} \quad (5.13)$$

For example, our hypothetical superintendent is contemplating reducing the student-teacher ratio by 2. Because the 95% confidence interval for  $\beta_1$  is  $[-3.30, -1.26]$ , the effect of reducing the student-teacher ratio by 2 could be as great as  $-3.30 \times (-2) = 6.60$  or as little as  $-1.26 \times (-2) = 2.52$ . Thus decreasing the student-teacher ratio by 2 is predicted to increase test scores by between 2.52 and 6.60 points, with a 95% confidence level.

## 5.3 Regression When $X$ Is a Binary Variable

The discussion so far has focused on the case that the regressor is a continuous variable. Regression analysis can also be used when the regressor is binary—that is, when it takes on only two values, 0 or 1. For example,  $X$  might be a worker's gender ( $= 1$  if female,  $= 0$  if male), whether a school district is urban or rural ( $= 1$  if urban,  $= 0$  if rural), or whether the district's class size is small or large ( $= 1$  if small,  $= 0$  if large). A binary variable is also called an **indicator variable** or sometimes a **dummy variable**.

### Interpretation of the Regression Coefficients

The mechanics of regression with a binary regressor are the same as if it is continuous. The interpretation of  $\beta_1$ , however, is different, and it turns out that regression with a binary variable is equivalent to performing a difference of means analysis, as described in Section 3.4.

To see this, suppose you have a variable  $D_i$  that equals either 0 or 1, depending on whether the student-teacher ratio is less than 20:

$$D_i = \begin{cases} 1 & \text{if the student-teacher ratio in } i^{\text{th}} \text{ district} < 20 \\ 0 & \text{if the student-teacher ratio in } i^{\text{th}} \text{ district} \geq 20 \end{cases} \quad (5.14)$$

The population regression model with  $D_i$  as the regressor is

$$Y_i = \beta_0 + \beta_1 D_i + u_i, \quad i = 1, \dots, n. \quad (5.15)$$