

# Chapter 15

## Inference with Regression Models

**LO 15.1:** Conduct tests of individual significance.

**LO 15.2:** Conduct a test of joint significance.

**LO 15.3:** Conduct a general test of linear restrictions.

**LO 15.4:** Calculate and interpret interval estimates for predictions.

**LO 15.5:** Explain the role of the assumptions on the OLS estimators.

**LO 15.6:** Describe common violations of the assumptions and offer remedies.

### 15.1 Tests of Significance

LO 15.1 Conduct tests of individual significance.

- With two explanatory variables to choose from, we can formulate three linear models:

$$\text{Model 1: } \text{Win} = \beta_0 + \beta_1 \text{BA} + \epsilon$$

$$\text{Model 2: } \text{Win} = \beta_0 + \beta_1 \text{ERA} + \epsilon$$

$$\text{Model 3: } \text{Win} = \beta_0 + \beta_1 \text{BA} + \beta_2 \text{ERA} + \epsilon$$

#### 15.1.1 Tests of Individual Significance

- Consider our standard multiple regression model: (??).
- In general, we can test whether  $\beta_j$  is equal to, greater than, or less than some hypothesized value  $\beta_{j0}$ .
- This test could have one of three forms:

Table 15.1: Three Forms of Individual Significance

Two-tailed Test	Right-tailed Test	Left-tailed Test
$H_0 : \beta_j = \beta_{j0}$	$H_0 : \beta_j \leq \beta_{j0}$	$H_0 : \beta_j \geq \beta_{j0}$
$H_A : \beta_j \neq \beta_{j0}$	$H_A : \beta_j > \beta_{j0}$	$H_A : \beta_j < \beta_{j0}$

### 15.1.2 The Test Statistic

- The appropriate test statistic is:

$$t_{df} = \frac{b_j - \beta_{j0}}{s_{b_j}} \quad (15.1)$$

- $s_{b_j}$  is the standard error of the estimator  $b_j$ .
- The test statistic will follow a  $t$ -distribution with degrees of freedom  $df = n - k - 1$ .

### 15.1.3 Testing $\beta_j = 0$

- By far, the most common hypothesis test for an individual coefficient is to test whether its value differs from zero.
- To see why, consider our model: (??).
- If a coefficient is equal to zero, then it implies that the explanatory variable is not a significant predictor of the response variable.

### 15.1.4 Computer-Generated Output

- Virtually all statistical software will automatically report a test statistic and a  $p$ -value with each coefficient estimate.
- These values can be used whether the regression coefficient differs from zero.
- To perform a one-sided test where the hypothesized value is zero, divide the computer-reported  $p$ -value in half.
- If we wish to test whether the coefficient differs from a non-zero value, we need to compute a new test statistic.

### 15.1.5 Intervals for the Parameters

- A confidence interval for the  $\beta_j$  parameter can be constructed using the formula:

$$b_j \pm t_{\alpha/2, df} s_{b_j} \quad (15.2)$$

- This can also be used to perform the two-sided test to determine whether a coefficient differs from zero.

- For ERA, the interval of  $[-0.15, -0.08]$  does not include 0, indicating ERA is a significant predictor.

### 15.1.6 Test for a Non-Zero Slope

- A capital asset pricing model follows the equation:

$$y = \alpha + \beta x + \epsilon \quad (15.3)$$

where  $y$  = the risk-adjusted return of an asset,  $R - R_f$  and  $x$  = the risk-adjusted return to the market  $R_M - R_f$ .

- The estimate of  $\beta$  is called the investment's beta value.
- A  $\beta > 1$  implies the stock is "aggressive", while a  $\beta < 1$  implies it is "conservative".

## 15.2 Test of Joint Significance

LO 15.2 Conduct a test of joint significance.

- In addition to conducting tests of individual significance, we may also want to test the joint significance of all  $k$  variables at once.
- The competing hypotheses for a test of joint significance are:

$$\begin{aligned} H_0 : \beta_1 &= \beta_2 = \cdots = \beta_k = 0 \\ H_A : \text{at least one } \beta_j &\neq 0 \end{aligned}$$

### 15.2.1 The Test Statistic

- The test statistic for a test of joint significant is

$$\begin{aligned} F(df1, df2) &= \frac{\frac{MSR}{MSE}}{\frac{SSR}{k}} \\ &= \frac{\frac{MSR}{MSE}}{\frac{SSE}{n-k-1}} \end{aligned} \quad (15.4)$$

where  $MSR$  and  $MSE$  are the mean regression sum of squares and the mean error sum of squares, respectively.

- The numerator degrees of freedom,  $df1 = k$ , while the denominator degrees of freedom,  $df2 = n - k - 1$ .

### 15.3 A General Test of Linear Restrictions

LO 15.3 Conduct a general test of linear restrictions.

- The significance tests in the previous section can also be labeled tests of **linear restrictions**.
- For example, if we have  $k = 3$  explanatory variables, testing whether  $\beta_2 = \beta_3 = 0$  is equivalent to testing whether to restrict the model to only  $x_1$ .
- In this section, we apply the  $F$ -test for any number of linear restrictions; the resulting  $F$ -test is often referred to as the **partial F**-test.

### 15.4 Interval Estimates for Predictions

LO 15.4 Calculate and interpret interval estimates for predictions.

### 15.5 Model Assumptions and Common Violations

LO 15.5 Explain the role of the assumptions on the OLS estimators.

- The statistical properties of OLS estimator, as well as the validity of the testing procedures, depend on a number of assumptions. We discuss those assumptions now.
  1. The model (??) is linear in the  $\beta$  parameters with an additive error  $\epsilon$ .
  2. Conditional on the  $x_1, \dots, x_k$  values, the expected error is 0, thus:

$$E(y) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k \quad (15.5)$$

- There is no exact linear relationship among the  $x_1, \dots, x_k$  values (no perfect multicollinearity).
- The variance of the error term  $\epsilon$  is the same for all  $x_1, \dots, x_k$  values. We call this homoskedasticity.
- The error term  $\epsilon$  is uncorrelated across observations, conditional on the explanatory variables. There is no serial correlation or autocorrelation.
- The error term  $\epsilon$  is not correlated with any of the predictors  $x_1, \dots, x_k$ . In other words, there is no endogeneity.
- The error term  $\epsilon$  is normally distributed. This assumption allows us to do hypothesis testing. If normality is not true, the tests may not be valid.

### 15.5.1 Checking the Assumptions

- The true error terms  $\epsilon$  cannot be observed because they exist only in the population. We can, however, look at the residuals,  $e = y - \hat{y}$ , where  $\hat{y} = b_0 + \sum b_i x_i$ , for each observation/
- It is common to plot residuals on the vertical axis and an explanatory variable on the horizontal axis.
- When estimating a regression in Excel, the dialog box that opens after choosing **Data > Data Analysis > Regression** allows us to select *Residuals* and Residual Plots options.

### 15.5.2 Common Violation 1: The Model Suffers from Multicollinearity

LO 15.6 Describe common violations of the assumptions and offer remedies.

- Perfect multicollinearity exists when two or more  $x$  variables have an exact linear relationship.
- For example, suppose the  $x$  data includes total cost, fixed cost and variable cost.
- Other data sets may have a great degree of multicollinearity that is not perfect.
- In these cases, we may see a high  $R^2$  coupled with individually insignificant explanatory variables. Additionally, unintuitive result may be indicative.
- A sample correlation between explanatory variables that is  $> 0.80$  or  $< -0.80$  suggests severe multicollinearity.

### 15.5.3 Remedyng Multicollinearity

- A good remedy may be simply drop one of the collinear variables if we can justify it as redundant.
- Alternatively, we could try to increase our sample size.
- Another option would be to try to transform our variables so that they are no longer collinear.
- Last, especially if we are interested only in maintaining a high predictive power, it may make sense to do nothing.

#### 15.5.4 Common Violation 2: The Error Term is Heteroskedastic

- The variance of the error term changes for different values of at least one explanatory variable.
- Informal residual plots can gauge heteroskedasticity (display a marked pattern).

#### 15.5.5 Remedyng Heteroskedastic

- Heteroskedastic results in inefficient estimators and the hypothesis tests for significance are no longer valid.
- To get around the second problem, some researchers use OLS estimates along with corrected standard errors, called White's standard errors. Many statistical packages have this option available, unfortunately the current version of Excel does not.

#### 15.5.6 Common Violation 3: The Error Term is Serially Correlated

- We assume that the error term is uncorrelated across observations when obtaining OLS estimates.
- But this often breaks down in time series data.
- Remedies are not easily accessible using Excel.

#### Common Violation 4: The Explanatory Variable is Endogenous

- Endogeneity in the regression model refers to the error term being correlated with the explanatory variables.
- This commonly occurs due to an omitted explanatory variable.
- For example, a person's salary may be highly correlated with that person's innate ability. But since we cannot include it, ability gets incorporated in the error term. If we try to predict salary by years of education, which may also be correlated with innate ability, then we have an endogeneity problem.
- Endogeneity will result in biased estimators, and so is quite a serious problem.
- Unfortunately, endogeneity is difficult to fix. Most commonly, we would like to find an instrumental variable, one that is correlated with the endogenous explanatory variable but uncorrelated with the error term. But it may be difficult to find such a variable.