

## Prelude

The simple and multiple linear regression methods we studied in Chapters 10 and 11 are used to model the relationship between a quantitative response variable and one or more explanatory variables. A key assumption for these models is that the deviations from the model fit are normally distributed. In this chapter we describe similar methods that are used when the response variable has only two possible values.

- How does the concentration of an insecticide relate to whether or not an insect is killed?
- To what extent does gender predict whether or not a college student will be a binge drinker?
- Is high blood pressure associated with an increased risk of death from cardiovascular disease?

CHAPTER 15

Logistic Regression

Our response variable has only two values: success or failure, live or die, acceptable or not. If we let the two values be 1 and 0, the mean is the proportion of ones,  $p = P(\text{success})$ . With  $n$  independent observations, we have the *binomial setting* (page 376). What is *new* here is that we have data on an *explanatory variable*  $x$ . We study how  $p$  depends on  $x$ . For example, suppose we are studying whether a patient lives ( $y = 1$ ) or dies ( $y = 0$ ) after being admitted to a hospital. Here,  $p$  is the probability that a patient lives, and possible explanatory variables include (a) whether the patient is in good condition or in poor condition, (b) the type of medical problem that the patient has, and (c) the age of the patient. Note that the explanatory variables can be either categorical or quantitative. Logistic regression<sup>1</sup> is a statistical method for describing these kinds of relationships.

## Binomial distributions and odds

In Chapter 5 we studied binomial distributions and in Chapter 8 we learned how to do statistical inference for the proportion  $p$  of successes in the binomial setting. We start with a brief review of some of these ideas that we will need in this chapter.

### EXAMPLE 15.1

Example 8.1 describes a survey of 17,096 students in U.S. four-year colleges. The researchers were interested in estimating the proportion of students who are frequent binge drinkers. A student who reports drinking five or more drinks in a row three or more times in the past two weeks is called a frequent binge drinker. In the notation of Chapter 5,  $p$  is the proportion of frequent binge drinkers in the entire population of college students in U.S. four-year colleges. The number of frequent binge drinkers in an SRS of size  $n$  has the binomial distribution with parameters  $n$  and  $p$ . The sample size is  $n = 17,096$  and the number of frequent binge drinkers in the sample is 3314. The sample proportion is

$$\hat{p} = \frac{3314}{17,096} = 0.1938$$

### odds

Logistic regressions work with **odds** rather than proportions. The odds are simply the ratio of the proportions for the two possible outcomes. If  $\hat{p}$  is the proportion for one outcome, then  $1 - \hat{p}$  is the proportion for the second outcome.

$$\text{ODDS} = \frac{\hat{p}}{1 - \hat{p}}$$

A similar formula for the population odds is obtained by substituting  $p$  for  $\hat{p}$  in this expression.

### EXAMPLE 15.2

For the binge-drinking data the proportion of frequent binge drinkers in the sample is  $\hat{p} = 0.1938$ , so the proportion of students who are not frequent binge drinkers is

$$1 - \hat{p} = 1 - 0.1938 = 0.8062$$

Therefore, the odds of a student being a frequent binge drinker are

$$\begin{aligned}\text{ODDS} &= \frac{\hat{p}}{1 - \hat{p}} \\ &= \frac{0.1938}{0.8062} \\ &= 0.24\end{aligned}$$

When people speak about odds, they often round to integers or fractions. Since 0.24 is approximately 1/4, we could say that the odds that a college student is a frequent binge drinker are 1 to 4. In a similar way, we could describe the odds that a college student is *not* a frequent binge drinker as 4 to 1.

In Example 8.8 (page 603) we compared the proportions of frequent binge drinkers among men and women college students using a confidence interval. There we found that the proportion for men was 0.227 (22.7%) and that the proportion for women was 0.170 (17.0%). The difference is 0.057 and the 95% confidence interval is (0.045, 0.069). We can summarize this result by saying, “The proportion of frequent binge drinkers is 5.7% higher among men than among women.”

Another way to analyze these data is to use logistic regression. The explanatory variable is gender, a categorical variable. To use this in a regression (logistic or otherwise), we need to use a numeric code. The usual way to do this is with an **indicator variable**. For our problem we will use an indicator of whether or not the student is a man:

indicator variable

$$x = \begin{cases} 1 & \text{if the student is a man} \\ 0 & \text{if the student is a woman} \end{cases}$$

The response variable is the proportion of frequent binge drinkers. For use in a logistic regression, we perform two transformations on this variable. First, we convert to odds. For men,

$$\begin{aligned}\text{ODDS} &= \frac{\hat{p}}{1 - \hat{p}} \\ &= \frac{0.227}{1 - 0.227} \\ &= 0.294\end{aligned}$$

Similarly, for women we have

$$\begin{aligned}\text{ODDS} &= \frac{\hat{p}}{1 - \hat{p}} \\ &= \frac{0.170}{1 - 0.170} \\ &= 0.205\end{aligned}$$

## The logistic regression model

In simple linear regression we modeled the mean  $\mu$  of the response variable  $y$  as a linear function of the explanatory variable:  $\mu = \beta_0 + \beta_1 x$ . With logistic

regression we are interested in modeling the mean of the response variable  $p$  in terms of an explanatory variable  $x$ . We could try to relate  $p$  and  $x$  through the equation  $p = \beta_0 + \beta_1 x$ . Unfortunately, this is not a good model. As long as  $\beta_1 \neq 0$ , extreme values of  $x$  will give values of  $\beta_0 + \beta_1 x$  that are inconsistent with the fact that  $0 \leq p \leq 1$ .

The logistic regression solution to this difficulty is to transform the odds ( $p/(1-p)$ ) using the natural logarithm. We use the term *log odds* for this transformation. We model the log odds as a linear function of the explanatory variable:

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x$$

Figure 15.1 graphs the relationship between  $p$  and  $x$  for some different values of  $\beta_0$  and  $\beta_1$ . For logistic regression we use *natural* logarithms. There are tables of natural logarithms, and many calculators have a built-in function for this transformation. As we did with linear regression, we use  $y$  for the response variable. So for men,

$$y = \log(\text{ODDS}) = \log(0.294) = -1.23$$

and for women,

$$y = \log(\text{ODDS}) = \log(0.205) = -1.59$$

In these expressions we use  $y$  as the observed value of the response variable, the log odds of being a frequent binge drinker. We are now ready to build the logistic regression model.

### Logistic Regression Model

The **statistical model for logistic regression** is

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x$$

where the  $p$  is a binomial proportion and  $x$  is the explanatory variable. The parameters of the logistic model are  $\beta_0$  and  $\beta_1$ .

#### EXAMPLE 15.3

For our binge-drinking example, there are  $n = 17,096$  students in the sample. The explanatory variable is gender, which we have coded using an indicator variable with values  $x = 1$  for men and  $x = 0$  for women. The response variable is also an indicator variable. Thus, the student is either a frequent binge drinker or the student is not a frequent binge drinker. Think of the process of randomly selecting a student and recording the values of  $x$  and whether or not the student is a frequent binge drinker. The model says that the probability ( $p$ ) that this student is a frequent binge drinker depends upon the student's gender ( $x = 1$  or  $x = 0$ ). So there are two possible values for  $p$ , say  $p_{\text{men}}$  and  $p_{\text{women}}$ .

Logistic regression with an indicator explanatory variable is a very special case. It is important because many multiple logistic regression analyses focus on one or more such variables as the primary explanatory variables of interest. For now, we use this special case to understand a little more about the model.

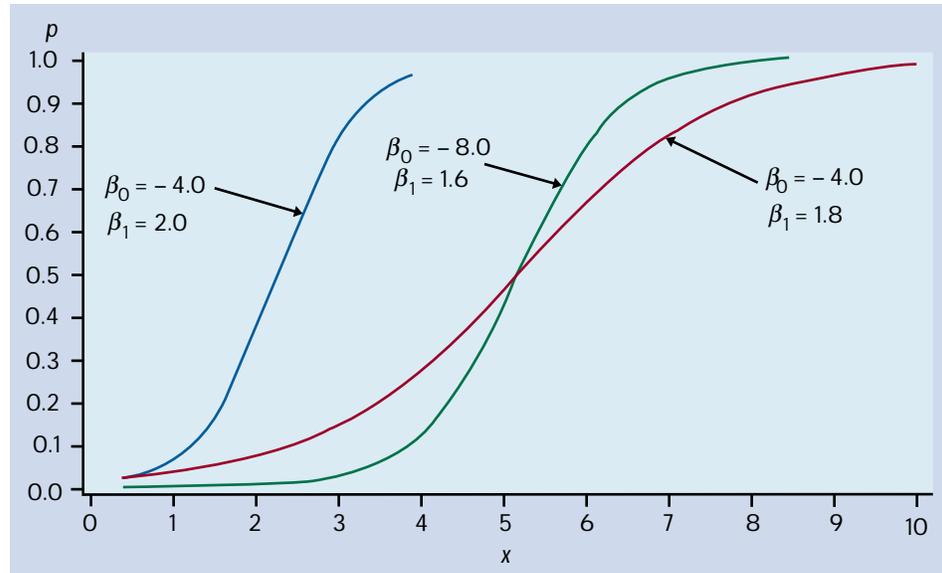


FIGURE 15.1 Plot of  $p$  versus  $x$  for selected values of  $\beta_0$  and  $\beta_1$ .

The logistic regression model specifies the relationship between  $p$  and  $x$ . Since there are only two values for  $x$ , we write both equations. For men,

$$\log\left(\frac{p_{\text{men}}}{1 - p_{\text{men}}}\right) = \beta_0 + \beta_1$$

and for women

$$\log\left(\frac{p_{\text{women}}}{1 - p_{\text{women}}}\right) = \beta_0$$

Note that there is a  $\beta_1$  term in the equation for men because  $x = 1$  but it is missing in the equation for women because  $x = 0$ .

### Fitting and interpreting the logistic regression model

In general the calculations needed to find estimates  $b_0$  and  $b_1$  for the parameters  $\beta_0$  and  $\beta_1$  are complex and require the use of software. When the explanatory variable has only two possible values, however, we can easily find the estimates. This simple framework also provides a setting where we can learn what the logistic regression parameters mean.

#### EXAMPLE 15.4

In the binge-drinking example, we found the log odds for men

$$y = \log\left(\frac{\hat{p}_{\text{men}}}{1 - \hat{p}_{\text{men}}}\right) = -1.23$$

and for women

$$y = \log\left(\frac{\hat{p}_{\text{women}}}{1 - \hat{p}_{\text{women}}}\right) = -1.59$$

The logistic regression model for men is

$$\log\left(\frac{p_{\text{men}}}{1 - p_{\text{men}}}\right) = \beta_0 + \beta_1$$

and for women, it is

$$\log\left(\frac{p_{\text{women}}}{1 - p_{\text{women}}}\right) = \beta_0$$

To find the estimates of  $b_0$  and  $b_1$ , we match the male and female model equations with the corresponding data equations. Thus, we see that the estimate of the intercept  $b_0$  is simply the log(ODDS) for the women:

$$b_0 = -1.59$$

and the slope is the difference between the log(ODDS) for the men and the log(ODDS) for the women:

$$b_1 = -1.23 - (-1.59) = 0.36$$

The fitted logistic regression model is

$$\log(\text{ODDS}) = -1.59 + 0.36x$$

The slope in this logistic regression model is the difference between the log(ODDS) for men and the log(ODDS) for women. Most people are not comfortable thinking in the log(ODDS) scale, so interpretation of the results in terms of the regression slope is difficult. Usually, we apply a transformation to help us. With a little algebra, it can be shown that

$$\frac{\text{ODDS}_{\text{men}}}{\text{ODDS}_{\text{women}}} = e^{0.36} = 1.43$$

odds ratio

The transformation  $e^{0.36}$  undoes the logarithm and transforms the logistic regression slope into an **odds ratio**, in this case, the ratio of the odds that a man is a frequent binge drinker to the odds that a woman is a frequent binge drinker. In other words, we can multiply the odds for women by the odds ratio to obtain the odds for men:

$$\text{ODDS}_{\text{men}} = 1.43 \times \text{ODDS}_{\text{women}}$$

In this case, the odds for men are 1.43 times the odds for women.

Notice that we have chosen the coding for the indicator variable so that the regression slope is positive. This will give an odds ratio that is greater than 1. Had we coded women as 1 and men as 0, the signs of the parameters would be reversed, the fitted equation would be  $\log(\text{ODDS}) = 1.59 - 0.36x$ , and the odds ratio would be  $e^{-0.36} = 0.70$ . The odds for women are 70% of the odds for men.

Logistic regression with an explanatory variable having two values is a very important special case. Here is an example where the explanatory variable is quantitative.

#### EXAMPLE 15.5

The CHEESE data set described in the Data Appendix includes a response variable called "Taste" that is a measure of the quality of the cheese obtained from several tasters. For this example, we will classify the cheese as acceptable (tasteok = 1) if

Taste  $\geq 37$  and unacceptable (tasteok = 0) if Taste  $< 37$ . This is our response variable. The data set contains three explanatory variables: “Acetic,” “H2S,” and “Lactic.” Let’s use Acetic as the explanatory variable. The model is

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x$$

where  $p$  is the probability that the cheese is acceptable and  $x$  is the value of Acetic. The model for estimated log odds fitted by software is

$$\log(\text{ODDS}) = b_0 + b_1 x = -13.71 + 2.25x$$

The odds ratio is  $e^{b_1} = 9.48$ . This means that if we increase the acetic acid content  $x$  by one unit, we increase the odds that the cheese will be acceptable by about 9.5 times.

## Inference for logistic regression

Statistical inference for logistic regression is very similar to statistical inference for simple linear regression. We calculate estimates of the model parameters and standard errors for these estimates. Confidence intervals are formed in the usual way, but we use standard normal  $z^*$ -values rather than critical values from the  $t$  distributions. The ratio of the estimate to the standard error is the basis for hypothesis tests. Often the test statistics are given as the squares of these ratios, and in this case the  $P$ -values are obtained from the chi-square distributions with 1 degree of freedom.

### Confidence Intervals and Significance Tests for Logistic Regression Parameters

A level  $C$  confidence interval for the slope  $\beta_1$  is

$$b_1 \pm z^* \text{SE}_{b_1}$$

The ratio of the odds for a value of the explanatory variable equal to  $x + 1$  to the odds for a value of the explanatory variable equal to  $x$  is the **odds ratio**.

A level  $C$  confidence interval for the odds ratio  $e^{\beta_1}$  is obtained by transforming the confidence interval for the slope:

$$(e^{b_1 - z^* \text{SE}_{b_1}}, e^{b_1 + z^* \text{SE}_{b_1}})$$

In these expressions  $z^*$  is the value for the standard normal density curve with area  $C$  between  $-z^*$  and  $z^*$ .

To test the hypothesis  $H_0: \beta_1 = 0$ , compute the **test statistic**

$$X^2 = \left(\frac{b_1}{\text{SE}_{b_1}}\right)^2$$

In terms of a random variable  $X^2$  having approximately a  $\chi^2$  distribution with 1 degree of freedom, the  $P$ -value for a test of  $H_0$  against  $H_a: \beta_1 \neq 0$  is  $P(\chi^2 \geq X^2)$ .

Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Square	Pr > Chi-Square	Odds Ratio
INTERCPT	1	-1.5869	0.0267	3520.4040	0.0001	.
X	1	0.3616	0.0388	86.6714	0.0001	1.436

FIGURE 15.2 Logistic regression output for the binge-drinking data, for Example 15.6.

We have expressed the hypothesis-testing framework in terms of the slope  $\beta_1$  because this form closely resembles what we studied in simple linear regression. In many applications, however, the results are expressed in terms of the odds ratio. A slope of 0 is the same as an odds ratio of 1, so we often express the null hypothesis of interest as “the odds ratio is 1.” This means that the two odds are equal and the explanatory variable is not useful for predicting the odds.

#### EXAMPLE 15.6

Figure 15.2 gives the output from the SAS logistic procedure for the binge-drinking example. The parameter estimates are given as  $b_0 = -1.5869$  and  $b_1 = 0.3616$ , the same as we calculated directly in Example 15.4, but with more significant digits. The standard errors are 0.0267 and 0.0388. A 95% confidence interval for the slope is

$$\begin{aligned} b_1 \pm z^*SE_{b_1} &= 0.3616 \pm (1.96)(0.0388) \\ &= 0.3616 \pm 0.0760 \end{aligned}$$

We are 95% confident that the slope is between 0.2855 and 0.4376. The output provides the odds ratio 1.436 but does not give the confidence interval. This is easy to compute from the interval for the slope:

$$\begin{aligned} (e^{b_1 - z^*SE_{b_1}}, e^{b_1 + z^*SE_{b_1}}) &= (e^{0.2855}, e^{0.4376}) \\ &= (1.33, 1.55) \end{aligned}$$

For this problem we would report, “College men are more likely to be frequent binge drinkers than college women (odds ratio = 1.44, 95% CI is 1.33 to 1.55).”

In applications such as these, it is standard to use 95% for the confidence coefficient. With this convention, the confidence interval gives us the result of testing the null hypothesis that the odds ratio is 1 for a significance level of 0.05. If the confidence interval does not include 1, we reject  $H_0$  and conclude that the odds for the two groups are different; if not, the data do not provide enough evidence to distinguish the groups in this way.

The following example is typical of many applications of logistic regression. Here there is a designed experiment with five different values for the explanatory variable.

#### EXAMPLE 15.7

An experiment was designed to examine how well the insecticide rotenone kills aphids that feed on the chrysanthemum plant called *Macrosiphoniella sanborni*.<sup>2</sup> The explanatory variable is the log concentration (in milligrams per liter) of the insecticide. At each concentration, approximately 50 insects were exposed. Each insect was either killed or not killed. We summarize the data using the number killed.

The response variable for logistic regression is the log odds of the proportion killed. Here are the data:

Concentration (log)	Number of insects	Number killed
0.96	50	6
1.33	48	16
1.63	46	24
2.04	49	42
2.32	50	44

If we transform the response variable (by taking log odds) and use least-squares, we get the fit illustrated in Figure 15.3. The logistic regression fit is given in Figure 15.4. It is a transformed version of Figure 15.3 with the fit calculated using the logistic model.

One of the major themes of this text is that we should present the results of a statistical analysis with a graph. For the insecticide example we have done this with Figure 15.4 and the results appear to be convincing. But suppose that rotenone has no ability to kill *Macrosiphoniella sanborni*. What is the chance that we would observe experimental results at least as convincing as what we observed if this supposition were true? The answer is the  $P$ -value for the test of the null hypothesis that the logistic regression slope is zero. If this  $P$ -value

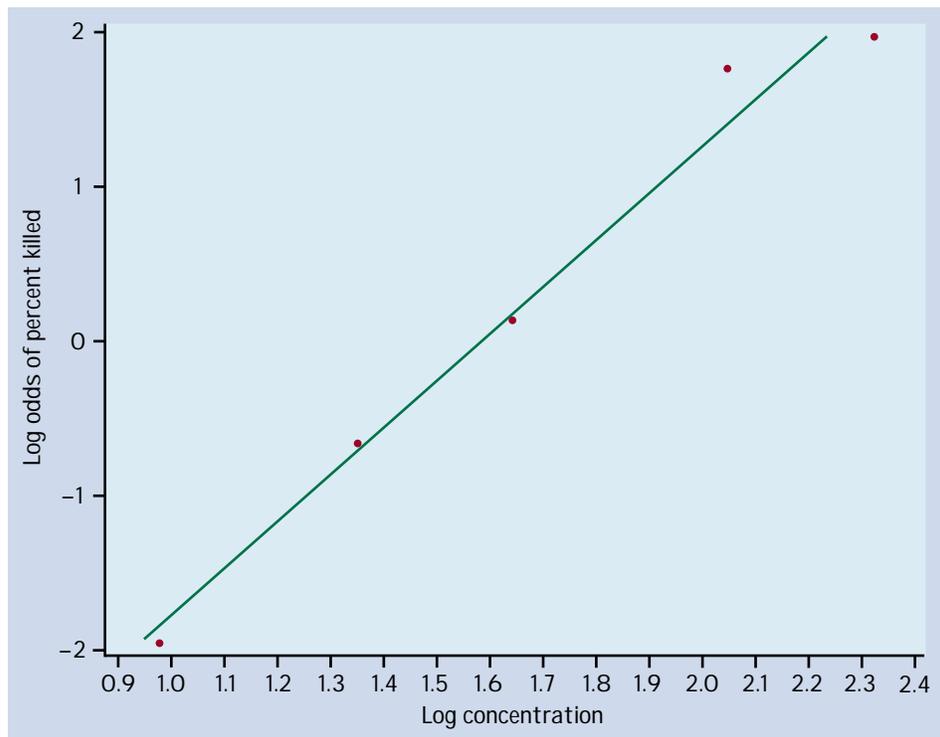


FIGURE 15.3 Plot of log odds of percent killed versus log concentration for the insecticide data, for Example 15.7.

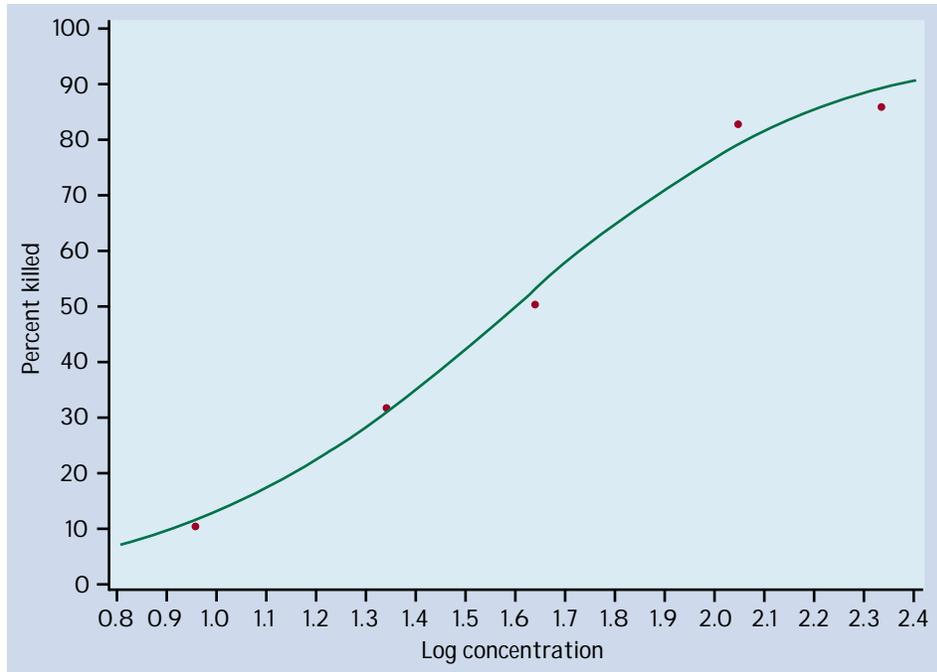


FIGURE 15.4 Plot of percent killed versus log concentration with the logistic fit for the insecticide data, for Example 15.7.

is not small, our graph may be misleading. Statistical inference provides what we need.

### EXAMPLE 15.8

The output produced by the SAS logistic procedure for the analysis of the insecticide data is given in Figure 15.5. The model is

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x$$

where the values of the explanatory variable  $x$  are 0.96, 1.33, 1.63, 2.04, 2.32. From the output we see that the fitted model is

$$\log(\text{ODDS}) = b_0 + b_1 x = -4.89 + 3.10x$$

This is the fit that we plotted in Figure 15.4. The null hypothesis that  $\beta_1 = 0$  is clearly rejected ( $X^2 = 64.07$ ,  $P < 0.001$ ). We calculate a 95% confidence interval for

Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Square	Pr > Chi-Square	Odds Ratio
INTERCPT	1	-4.8869	0.6429	57.7757	0.0001	.
LCONC	1	3.1035	0.3877	64.0744	0.0001	22.277

FIGURE 15.5 Logistic regression output for the insecticide data, for Example 15.8.

$\beta_1$  using the estimate  $b_1 = 3.1035$  and its standard error  $SE_{b_1} = 0.3877$  given in the output:

$$\begin{aligned} b_1 \pm z^* SE_{b_1} &= 3.1035 \pm (1.96)(0.3877) \\ &= 3.1035 \pm 0.7599 \end{aligned}$$

We are 95% confident that the true value of the slope is between 2.34 and 3.86.

The odds ratio is given on the output as 22.277. An increase of one unit in the log concentration of insecticide ( $x$ ) is associated with a 22-fold increase in the odds that an insect will be killed. The confidence interval for the odds is obtained from the interval for the slope:

$$\begin{aligned} (e^{b_1 - z^* SE_{b_1}}, e^{b_1 + z^* SE_{b_1}}) &= (e^{2.34361}, e^{3.86339}) \\ &= (10.42, 47.63) \end{aligned}$$

Note again that the test of the null hypothesis that the slope is zero is the same as the test of the null hypothesis that the odds are 1. If we were reporting the results in terms of the odds, we could say, "The odds of killing an insect increase by a factor of 22.3 for each unit increase in the log concentration of insecticide ( $X^2 = 64.07$ ,  $P < 0.001$ ; 95% CI is 10.4 to 47.6)."

In Example 15.5 we studied the problem of predicting whether or not the taste of cheese was acceptable using Acetic as the explanatory variable. We now revisit this example and show how statistical inference is an important part of the conclusion.

### EXAMPLE 15.9

The output for a logistic regression analysis using Acetic as the explanatory variable is given in Figure 15.6. In Example 15.5 we gave the fitted model:

$$\log(\text{ODDS}) = b_0 + b_1 x = -13.71 + 2.25x$$

From the output we see that because  $P = 0.0285$ , we can reject the null hypothesis that  $\beta_1 = 0$ . The value of the test statistic is  $X^2 = 4.79$  with 1 degree of freedom. We use the estimate  $b_1 = 2.2490$  and its standard error  $SE_{b_1} = 1.0271$  to compute the 95% confidence interval for  $\beta_1$ :

$$\begin{aligned} b_1 \pm z^* SE_{b_1} &= 2.2490 \pm (1.96)(1.0271) \\ &= 2.2490 \pm 2.0131 \end{aligned}$$

Our estimate of the slope is 2.25 and we are 95% confident that the true value is between 0.24 and 4.26. For the odds ratio, the estimate on the output is 9.48. The 95% confidence interval is

$$\begin{aligned} (e^{b_1 - z^* SE_{b_1}}, e^{b_1 + z^* SE_{b_1}}) &= (e^{0.23588}, e^{4.26212}) \\ &= (1.27, 70.96) \end{aligned}$$

Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Square	Pr > Chi-Square	Odds Ratio
INTERCPT	1	-13.7052	5.9319	5.3380	0.0209	.
ACETIC	1	2.2490	1.0271	4.7947	0.0285	9.479

FIGURE 15.6 Logistic regression output for the cheese data with Acetic as the explanatory variable, for Example 15.9.

We estimate that increasing the acetic acid content of the cheese by one unit will increase the odds that the cheese will be acceptable by about 9 times. The data, however, do not give us a very accurate estimate. The odds ratio could be as small as a little more than 1 or as large as 71 with 95% confidence. We have evidence to conclude that cheeses with higher concentrations of acetic acid are more likely to be acceptable, but establishing the true relationship accurately would require more data.

## Multiple logistic regression

The cheese example that we just considered naturally leads us to the next topic. The data set includes three variables: Acetic, H2S, and Lactic. We examined the model where Acetic was used to predict the odds that the cheese was acceptable. Do the other explanatory variables contain additional information that will give us a better prediction? We use **multiple logistic regression** to answer this question. Generating the computer output is easy, just as it was when we generalized simple linear regression with one explanatory variable to multiple linear regression with more than one explanatory variable in Chapter 11. The statistical concepts are similar although the computations are more complex. Here is the example.

multiple logistic regression

### EXAMPLE 15.10

As in Example 15.8, we predict the odds that the cheese is acceptable. The explanatory variables are Acetic, H2S, and Lactic. Figure 15.7 gives the output. The fitted model is

$$\begin{aligned}\log(\text{ODDS}) &= b_0 + b_1\text{Acetic} + b_2\text{H2S} + b_3\text{Lactic} \\ &= -14.26 + 0.58\text{Acetic} + 0.68\text{H2S} + 3.47\text{Lactic}\end{aligned}$$

When analyzing data using multiple regression, we first examine the hypothesis that all of the regression coefficients for the explanatory variables are zero. We do the same for logistic regression. The hypothesis

$$H_0: \beta_1 = \beta_2 = \beta_3 = 0$$

Criterion	Intercept Only	Intercept and Covariates	Chi-Square for Covariates			
-2 LOG L	34.795	18.461	16.334 with 3 DF (p=0.0010)			
Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Square	Pr > Chi-Square	Odds Ratio
INTERCPT	1	-14.2604	8.2869	2.9613	0.0853	.
ACETIC	1	0.5845	1.5442	0.1433	0.7051	1.794
H2S	1	0.6849	0.4040	2.8730	0.0901	1.983
LACTIC	1	3.4684	2.6497	1.7135	0.1905	32.086

FIGURE 15.7 Logistic regression output for the cheese data with Acetic, H2S, and Lactic as the explanatory variables, for Example 15.10.

is tested by a chi-square statistic with 3 degrees of freedom. This is given in the output on the line for the criterion “ $-2 \text{ LOG } L$ ” under the heading “Chi-Square for Covariates.” The statistic is  $X^2 = 16.33$  and the  $P$ -value is 0.001. We reject  $H_0$  and conclude that one or more of the explanatory variables can be used to predict the odds that the cheese is acceptable. We now examine the coefficients for each variable and the tests that each of these is 0. The  $P$ -values are 0.71, 0.09, and 0.19. None of the null hypotheses,  $H_0: \beta_1 = 0$ ,  $H_0: \beta_2 = 0$ , and  $H_0: \beta_3 = 0$ , can be rejected.

Our initial multiple logistic regression analysis told us that the explanatory variables contain information that is useful for predicting whether or not the cheese is acceptable. Because the explanatory variables are correlated, however, we cannot clearly distinguish which variables or combinations of variables are important. Further analysis of these data using subsets of the three explanatory variables is needed to clarify the situation. We leave this work for the exercises.

## SUMMARY

If  $\hat{p}$  is the sample proportion, then the **odds** are  $\hat{p}/(1 - \hat{p})$ , the ratio of the proportion of times the event happens to the proportion of times the event does not happen.

The **logistic regression model** relates the log of the odds to the explanatory variable:

$$\log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 x_i$$

where the response variables for  $i = 1, 2, \dots, n$  are independent binomial random variables with parameters 1 and  $p_i$ ; that is, they are independent with distributions  $B(1, p_i)$ . The explanatory variable is  $x$ .

The **parameters** of the logistic model are  $\beta_0$  and  $\beta_1$ .

The **odds ratio** is  $e^{\beta_1}$ , where  $\beta_1$  is the slope in the logistic regression model.

A **level  $C$  confidence interval for the intercept**  $\beta_0$  is

$$b_0 \pm z^* \text{SE}_{b_0}$$

A **level  $C$  confidence interval for the slope**  $\beta_1$  is

$$b_1 \pm z^* \text{SE}_{b_1}$$

A **level  $C$  confidence interval for the odds ratio**  $e^{\beta_1}$  is obtained by transforming the confidence interval for the slope:

$$(e^{b_1 - z^* \text{SE}_{b_1}}, e^{b_1 + z^* \text{SE}_{b_1}})$$

In these expressions  $z^*$  is the value for the standard normal density curve with area  $C$  between  $-z^*$  and  $z^*$ .

To test the hypothesis  $H_0: \beta_1 = 0$ , compute the **test statistic**

$$\chi^2 = \left( \frac{b_1}{\text{SE}_{b_1}} \right)^2$$

In terms of a random variable  $X^2$  having a  $\chi^2$  distribution with 1 degree of freedom, the  $P$ -value for a test of  $H_0$  against  $H_a: \beta_1 \neq 0$  is  $P(\chi^2 \geq X^2)$ . This is the same as testing the null hypothesis that the odds ratio is 1.

In **multiple logistic regression** the response variable has two possible values, as in logistic regression, but there can be several explanatory variables.

## CHAPTER 15 EXERCISES

- 15.1** There is much evidence that high blood pressure is associated with increased risk of death from cardiovascular disease. A major study of this association examined 3338 men with high blood pressure and 2676 men with low blood pressure. During the period of the study, 21 men in the low-blood-pressure and 55 in the high-blood-pressure group died from cardiovascular disease.
- Find the proportion of men who died from cardiovascular disease in the high-blood-pressure group. Then calculate the odds.
  - Do the same for the low-blood-pressure group.
  - Now calculate the odds ratio with the odds for the high-blood-pressure group in the numerator. Describe the result in words.
- 15.2** To what extent do syntax textbooks, which analyze the structure of sentences, illustrate gender bias? A study of this question sampled sentences from ten texts. One part of the study examined the use of the words “girl,” “boy,” “man,” and “woman.” We will call the first two words juvenile and the last two adult. Here are data from one of the texts. (From Monica Macaulay and Colleen Brice, “Don’t touch my projectile: gender bias and stereotyping in syntactic examples,” *Language*, 73, no. 4 (1997), pp. 798–825.)

Gender	$n$	$X$ (juvenile)
Female	60	48
Male	132	52

- Find the proportion of the female references that are juvenile. Then transform this proportion to odds.
- Do the same for the male references.
- What is the odds ratio for comparing the female references to the male references? (Put the female odds in the numerator.)

- 15.3** Refer to the study of cardiovascular disease and blood pressure in Exercise 15.1. Computer output for a logistic regression analysis of these data gives the estimated slope  $b_1 = 0.7505$  with standard error  $SE_{b_1} = 0.2578$ .
- Give a 95% confidence interval for the slope.
  - Calculate the  $X^2$  statistic for testing the null hypothesis that the slope is zero and use Table F to find an approximate  $P$ -value.
  - Write a short summary of the results and conclusions.
- 15.4** The data from the study of gender bias in syntax textbooks given in Exercise 15.2 are analyzed using logistic regression. The estimated slope is  $b_1 = 1.8171$  and its standard error is  $SE_{b_1} = 0.3686$ .
- Give a 95% confidence interval for the slope.
  - Calculate the  $X^2$  statistic for testing the null hypothesis that the slope is zero and use Table F to find an approximate  $P$ -value.
  - Write a short summary of the results and conclusions.
- 15.5** The results describing the relationship between blood pressure and cardiovascular disease are given in terms of the change in log odds in Exercise 15.3.
- Transform the slope to the odds and the 95% confidence interval for the slope to a 95% confidence interval for the odds.
  - Write a conclusion using the odds to describe the results.
- 15.6** The gender bias in syntax textbooks is described in the log odds scale in Exercise 15.4.
- Transform the slope to the odds and the 95% confidence interval for the slope to a 95% confidence interval for the odds.
  - Write a conclusion using the odds to describe the results.
- 15.7** To be competitive in global markets, many U.S. corporations are undertaking major reorganizations. Often these involve “downsizing” or a “reduction in force” (RIF), where substantial numbers of employees are terminated. Federal and various state laws require that employees be treated equally regardless of their age. In particular, employees over the age of 40 years are in a “protected” class, and many allegations of discrimination focus on comparing employees over 40 with their younger coworkers. Here are the data for a recent RIF:

	Over 40	
Terminated	No	Yes
Yes	7	41
No	504	765

- (a) Write the logistic regression model for this problem using the log odds of a RIF as the response variable and an indicator for over and under 40 years of age as the explanatory variable.
- (b) Explain the assumptions concerning binomial distributions in terms of the variables in this exercise. To what extent do you think that these assumptions are reasonable?
- (c) Software gives the estimated slope  $b_1 = 1.3504$  and its standard error  $SE_{b_1} = 0.4130$ . Transform the results to the odds scale. Summarize the results and write a short conclusion.
- (d) If additional explanatory variables were available, for example, a performance evaluation, how would you use this information to study the RIF?

- 15.8** A study of alcohol use and deaths due to bicycle accidents collected data on a large number of fatal accidents. For each of these, the individual who died was classified according to whether or not there was a positive test for alcohol and by gender. Here are the data. (From Guohua Li and Susan P. Baker, "Alcohol in fatally injured bicyclists," *Accident Analysis and Prevention*, 26 (1994), pp. 543–548.)

Gender	$n$	$X$ (tested positive)
Female	191	27
Male	1520	515

Use logistic regression to study the question of whether or not gender is related to alcohol use in people who are fatally injured in bicycle accidents.

- 15.9** In Examples 15.5 and 15.9, we analyzed data from the CHEESE data set described in the Data Appendix. In those examples, we used Acetic as the explanatory variable. Run the same analysis using H2S as the explanatory variable.
- 15.10** Refer to the previous exercise. Run the same analysis using Lactic as the explanatory variable.
- 15.11** For the cheese data analyzed in Examples 15.5, 15.9, and 15.10 and in the two exercises above, there are three explanatory variables. There are three different logistic regressions that include two explanatory variables. Run these. Summarize the results of these analyses, the ones using each explanatory variable alone, and all three explanatory variables together. What do you conclude?

*The following four exercises use the CSDATA data set described in the Data Appendix. We examine models for relating success as measured by the GPA to several explanatory variables. In Chapter 11 we used multiple regression methods for our analysis. Here, we define an indicator variable, say HIGPA, to be 1 if the GPA is 3.0 or better and 0 otherwise.*

- 15.12** Use a logistic regression to predict HIGPA using the three high school grade summaries as explanatory variables.
- Summarize the results of the hypothesis test that the coefficients for all three explanatory variables are zero.
  - Give the coefficient for high school math grades with a 95% confidence interval. Do the same for the two other predictors in this model.
  - Summarize your conclusions based on parts (a) and (b).
- 15.13** Use a logistic regression to predict HIGPA using the two SAT scores as explanatory variables.
- Summarize the results of the hypothesis test that the coefficients for both explanatory variables are zero.
  - Give the coefficient for SAT math with a 95% confidence interval. Do the same for the SAT verbal score.
  - Summarize your conclusions based on parts (a) and (b).
- 15.14** Run a logistic regression to predict HIGPA using the three high school grade summaries and the two SAT scores as explanatory variables. We want to produce an analysis that is similar to that done for the case study in Chapter 11.
- Test the null hypothesis that the coefficients of the three high school grade summaries are zero; that is, test  $H_0: \beta_{\text{HSM}} = \beta_{\text{HSS}} = \beta_{\text{HSE}} = 0$ .
  - Test the null hypothesis that the coefficients of the two SAT scores are zero; that is, test  $H_0: \beta_{\text{SATM}} = \beta_{\text{SATV}} = 0$ .
  - What do you conclude from the tests in (a) and (b)?
- 15.15** In this exercise we investigate the effect of gender on the odds of getting a high GPA.
- Use gender to predict HIGPA using a logistic regression. Summarize the results.
  - Perform a logistic regression using gender and the two SAT scores to predict HIGPA. Summarize the results.
  - Compare the results of parts (a) and (b) with respect to how gender relates to HIGPA. Summarize your conclusions.
- 15.16** In Example 2.32 (page 189) we studied an example of Simpson's paradox, *the reversal of the direction of a comparison or an association when data from several groups are combined to form a single group*. The data concerned two hospitals, A and B, and whether or not patients undergoing surgery died or survived. Here are the data for all patients:

	Hospital A	Hospital B
Died	63	16
Survived	2037	784
Total	2100	800

And here are the more detailed data where the patients are categorized as being in good condition or poor condition:

	Good condition		Poor condition	
	Hospital A	Hospital B	Hospital A	Hospital B
Died	6	8	57	8
Survived	594	592	1443	192
Total	600	600	1500	200

- (a) Use a logistic regression to model the odds of death with hospital as the explanatory variable. Summarize the results of your analysis and give a 95% confidence interval for the odds ratio of Hospital A relative to Hospital B.
- (b) Rerun your analysis in (a) using hospital and the condition of the patient as explanatory variables. Summarize the results of your analysis and give a 95% confidence interval for the odds ratio of Hospital A relative to Hospital B.
- (c) Explain Simpson's paradox in terms of your results in parts (a) and (b).

## NOTES

1. Logistic regression models for the general case where there are more than two possible values for the response variable have been developed. These are considerably more complicated and are beyond the scope of our present study. For more information on logistic regression, see A. Agresti, *An Introduction to Categorical Data Analysis*, Wiley, New York, 1996; and D. W. Hosmer and S. Lemeshow, *Applied Logistic Regression*, Wiley, New York, 1989.
2. This example is taken from a classic text written by a contemporary of R. A. Fisher, the person who developed many of the fundamental ideas of statistical inference that we use today. The reference is D. J. Finney, *Probit Analysis*, Cambridge University Press, Cambridge, 1947. Although not included in our analysis, it is important to note that the experiment included a control group that received no insecticide. No aphids died in this group. We have chosen to call the response "dead." In Finney's text, the category is described as "apparently dead, moribund, or so badly affected as to be unable to walk more than a few steps." This is an early example of the need to make careful judgments when defining variables to be used in a statistical analysis. An insect that is "unable to walk more than a few steps" is unlikely to eat very much of a chrysanthemum plant!