



GoodLifeStudio/Getty Images

Logistic Regression

Introduction

The linear regression methods we studied in Chapters 12 and 13 are used to model the relationship between a quantitative response variable and one or more explanatory variables. In this chapter, we describe similar methods for use when the response variable *has only two possible outcomes*. For example,

- BeyondtheRack.com and HauteLook.com are two online destinations offering flash sale events. A response variable of interest to each of their sales divisions is whether a customer does or does not buy the flash sale item.
- For the JPMorgan Chase & Co. Human Resources Division, a response variable of interest is whether a candidate accepts or declines a job offer.

In general, we call the two outcomes of the response variable “success” and “failure” and represent them by 1 (for a success) and 0 (for a failure). The mean is then the proportion of 1s, $p = P(\text{success})$.

If our data are n independent observations with the same p , we are in the *binomial setting* and can use the inference methods of Chapter 10. What is *new* in this chapter is that the data now include at least one *explanatory variable* x and the probability p of a success depends on the value of x . The explanatory variables can be either categorical or quantitative. For example, the probability a customer purchases the flash sale item could depend on the age and gender of the customer, as well as the type of clothing item on sale and the percent discount. The probability a candidate accepts a job offer from JPMorgan Chase & Co. could depend on the salary amount, the level of guaranteed bonuses, and whether the offer includes a noncompete clause.

CHAPTER OUTLINE

- 18.1 The Logistic Regression Model
- 18.2 Inference for Logistic Regression
- 18.3 Multiple Logistic Regression



binomial setting, p. 250

Because it is now a probability that depends on explanatory variables, inference methods are needed to ensure that the probability p satisfies $0 \leq p \leq 1$. Logistic regression is a statistical method for describing these kinds of relationships.¹ Just as we did with linear regression, we start our study considering a single explanatory variable (simple logistic regression) and then expand to consider many variables (multiple logistic regression).

18.1 The Logistic Regression Model

When you complete this section, you will be able to:

- Find the odds from a single probability.
- Describe the statistical model for logistic regression with a single explanatory variable.
- Use the log odds to obtain regression parameter estimates in cases involving an indicator explanatory variable.
- Find the odds ratio for comparing two proportions.

In general, the data for simple logistic regression are n independent cases, each consisting of a value of the explanatory variable x and either a success or a failure for that trial. For example, x may be the offered salary amount, and “success” means that this applicant accepted the job offer. Every observation may have a different value of x .

To introduce logistic regression, however, it is convenient to start with the special case in which the explanatory variable x is also a two outcome variable. The data then contain a number of outcomes (success or failure) for each of the two values of x . There are also just two values of p , one for each value of x . Assuming the count of successes for each value of x has a binomial distribution, we are on familiar ground as described in Chapters 5 and 10. Here is an example.



Neiron Photo/Shutterstock



Clothing Color and Tipping What are the factors that affect a customer's tipping behavior? Studies have shown that a server's gender and various aspects of a server's appearance have an effect on tipping, unrelated to the quality of service. Some of these same studies have shown that the effect is different for male and female customers.

Because the color red has been shown to increase the physical attractiveness of women, a group of researchers decided to see if the color of clothing a female server wears has an effect on the tipping behavior.² Although they considered both male and female customers, we focus on the 418 male customers in their study.

The response variable is whether the male customer left a tip. The explanatory variable is whether the female server wore a red top. Let's express this condition numerically using an indicator variable,

$$x = \begin{cases} 1 & \text{if the server wore a red top} \\ 0 & \text{if the server wore a different colored top} \end{cases}$$

The female servers in the study wore a red top for 69 of the customers and wore a different colored top for the other 349.

The probability that a randomly chosen customer will tip has two values, p_1 for those whose server wore a red top and p_0 for those whose server wore a different colored top. The number of customers who tipped a server wearing a red top has the binomial distribution $B(69, p_1)$. The number of customers who tipped a server wearing a different color top has the $B(349, p_0)$ distribution. ■

Binomial distributions and odds

We begin with a review of some ideas associated with binomial distributions.

EXAMPLE 18.1



RED

CASE 18.1 **Proportion of Tippers** In Chapter 10, we used sample proportions to estimate population proportions. For this study, 40 of the 69 male customers tipped a server who was wearing red and 130 of the 349 customers tipped a server who was wearing a different color. The following table summarizes these results.

	Observed numbers of male customers		
	Server wearing red		Total
Tipped the server	No	Yes	
No	219	29	248
Yes	130	40	170
Total	349	69	418

Our estimates of the two population proportions are

$$\text{red: } \hat{p}_1 = \frac{40}{69} = 0.5797$$

and

$$\text{not red: } \hat{p}_0 = \frac{130}{349} = 0.3725$$

That is, we estimate that 58.0% of the male customers will tip if the server wears red, and 37.3% of the male customers will tip if the server wears a different color. ■

odds Logistic regression works with odds rather than proportions. The **odds** are the ratio of the proportions for the two possible outcomes. If p is the probability of a success, then $1 - p$ is the probability of a failure, and

$$\text{odds} = \frac{p}{1 - p} = \frac{\text{probability of success}}{\text{probability of failure}}$$

A similar formula for the sample odds is obtained by substituting \hat{p} for p in this expression. Let's now compute these odds for our tipping study.

EXAMPLE 18.2

CASE 18.1 **Odds of Tipping** The proportion of tippers among male customers who have a server wearing red is $\hat{p}_1 = 0.5797$, so the proportion of male customers who are not tippers when their server wears red is

$$1 - \hat{p}_1 = 1 - 0.5797 = 0.4203$$

The sample odds of a male customer tipping when the server wears red are, therefore,

$$\begin{aligned} \text{odds} &= \frac{\hat{p}_1}{1 - \hat{p}_1} \\ &= \frac{0.5797}{1 - 0.5797} = 1.3793 \end{aligned}$$

For the case when the server does not wear red, the sample odds are

$$\begin{aligned}\text{odds} &= \frac{\hat{p}_0}{1 - \hat{p}_0} \\ &= \frac{0.3725}{1 - 0.3725} = 0.5936 \blacksquare\end{aligned}$$

When people speak about odds, they often round to integers or fractions. Because 1.3793 is approximately 7/5, we could say that the odds that a male customer tips when the server wears red are 7 to 5, or for every 7 males who do tip, there are 5 males who do not. In a similar way, we could describe the odds that a male customer does *not* tip when the server wears red as 5 to 7. We could also say the odds that a male customer tips when the server does not wear red is about 6 to 10.

APPLY YOUR KNOWLEDGE

18.1 Energy drink commercials. A study was designed to compare Monster Energy TV commercials. Each participant was shown the commercials, A and B, in random order and asked to select the better one. There were 150 women and 220 men who participated in the study. Commercial A was selected by 90 women and by 97 men. Find the odds of selecting Commercial A for the women. Do the same for the men.

18.2 Use of audio/visual sharing through social media. In Case 10.3 (page 525), we studied data on large and small food and beverage companies and the use of audio/visual sharing through social media. Here are the data:

Observed numbers of companies			
Use A/V sharing	Size		Total
	Small	Large	
No	28	25	53
Yes	150	27	177
Total	178	52	230

What proportion of the small companies use audio/visual sharing?
 What proportion of the large companies use audio/visual sharing?
 Convert each of these proportions to odds.

Model for logistic regression

In Section 10.2 (page 523), we learned how to compare the proportions of two groups (such as large and small companies) using z tests and confidence intervals. Simple logistic regression is another way to make this comparison, but it extends to more general settings when working with a success-or-failure response variable.

In simple linear regression we modeled the mean μ of the response variable y as a linear function of the explanatory variable: $\mu = \beta_0 + \beta_1 x$. When y is just 1 or 0 (success or failure), the mean is the probability p of a success. Simple logistic regression models the mean p in terms of an explanatory variable x . We might try to relate p and x as in simple linear regression: $p = \beta_0 + \beta_1 x$. Unfortunately, this is not a good model. Whenever $\beta_1 \neq 0$, extreme values of x will give values of $\beta_0 + \beta_1 x$ that fall outside the range of possible values of p , $0 \leq p \leq 1$.

The logistic regression model removes this difficulty by working with the natural logarithm of the odds. We use the term **log odds** or **logit** for

this transformation of p . We model the log odds as a linear function of the explanatory variable:

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

As p moves from 0 to 1, the log odds move through all negative and positive numerical values. The logit of $p = 0.5$ is 0 so negative logits refer to $p < 0.5$ and positive logits refer to $p > 0.5$. Here is a summary of the simple logistic regression model.

SIMPLE LOGISTIC REGRESSION MODEL

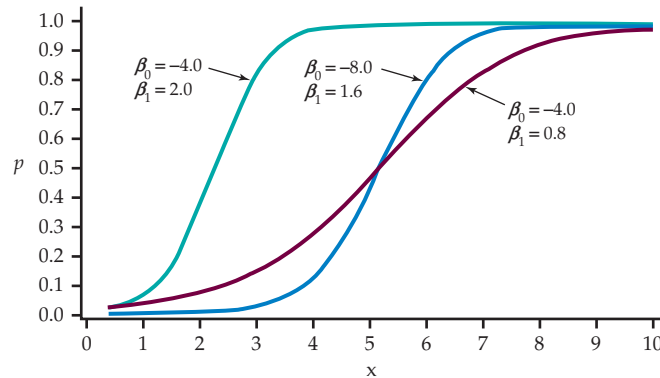
The **statistical model for simple logistic regression** is

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

where p is a binomial proportion and x is the explanatory variable. The parameters of the logistic model are β_0 and β_1 .

Figure 18.1 graphs the relationship between p and x for some different values of β_0 and β_1 . The logistic regression model uses *natural* logarithms. Most spreadsheets and calculators have a built-in function for the natural logarithm, often labeled “ln.”

FIGURE 18.1 Plot of p versus x for selected values of β_0 and β_1 .



Returning to the tipping study, for servers wearing red we have

$$\log(\text{odds}) = \log(1.3793) = 0.3216$$

and for servers not wearing red we have

$$\log(\text{odds}) = \log(0.5936) = -0.5215$$

Verify these results, remembering that “log” in these equations is the natural logarithm.

APPLY YOUR KNOWLEDGE

18.3 Log odds choosing Commercial A. Refer to Exercise 18.1. Find the log odds for the women and the log odds for the men choosing Commercial A.

18.4 Log odds for use of audio/visual sharing. Refer to Exercise 18.2. Find the log odds for the small and large companies.

Fitting and interpreting the logistic regression model

We must now fit the logistic regression model to data. In general, the data consist of n observations on the explanatory variable x , each with a success-or-failure response. Our tipping example has an indicator (0 or 1) explanatory variable. Logistic regression with an indicator explanatory variable is a special case but is important in practice. We use this special case to understand a little more about the model.

EXAMPLE 18.3



Logistic Model for Tipping Behavior In the tipping example, there are $n = 418$ observations. The explanatory variable is whether the server wore a red top, which we coded using an indicator variable with values $x = 1$ for servers wearing red and $x = 0$ for servers wearing a different color. There are 69 observations with $x = 1$ and 349 observations with $x = 0$. The response variable is also an indicator variable: $y = 1$ if the customer left a tip and $y = 0$ if not. The model says that the probability p of leaving a tip depends on the color of the server's top ($x = 1$ or $x = 0$). There are two possible values for p —say, p_1 for servers wearing red and p_0 for servers wearing a different color. The model says that for servers wearing red

$$\log\left(\frac{p_1}{1-p_1}\right) = \beta_0 + \beta_1(1) = \beta_0 + \beta_1$$

and for servers wearing a different color

$$\log\left(\frac{p_0}{1-p_0}\right) = \beta_0 + \beta_1(0) = \beta_0$$

There is a β_1 term in the equation for servers wearing red because $x = 1$. It is missing in the equation for servers wearing a different color because $x = 0$. ■

In general, the calculations needed to find the estimates b_0 and b_1 for the parameters β_0 and β_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 18.4



Parameter Estimates for Tipping Behavior For the tipping example, we found the log odds for servers wearing red,

$$\log\left(\frac{\hat{p}_1}{1-\hat{p}_1}\right) = 0.3216$$

and for servers wearing a different color,

$$\log\left(\frac{\hat{p}_0}{1-\hat{p}_0}\right) = -0.5215$$

To find estimates b_0 and b_1 of the model parameters β_0 and β_1 , we match the two model equations in Example 18.3 with the corresponding data equations. Because

$$\log\left(\frac{p_0}{1-p_0}\right) = \beta_0 \quad \text{and} \quad \log\left(\frac{\hat{p}_0}{1-\hat{p}_0}\right) = -0.5215$$

the estimate b_0 of the intercept is simply the $\log(\text{odds})$ for servers wearing a different color,

$$b_0 = -0.5215$$

Similarly, the estimated slope is the difference between the $\log(\text{odds})$ for servers wearing red and the $\log(\text{odds})$ for servers wearing a different color,

$$b_1 = 0.3216 - (-0.5215) = 0.8431$$

The fitted logistic regression model is

$$\log(\text{odds}) = -0.5215 + 0.8431x \blacksquare$$

The slope in this logistic regression model is the difference between the $\log(\text{odds})$ for servers wearing red and the $\log(\text{odds})$ for servers wearing a different color. Most people are not comfortable thinking in the $\log(\text{odds})$ scale, so interpretation of the results in terms of the regression slope is difficult.

EXAMPLE 18.5



Transforming Estimates to the Odds Scale To get to the odds scale, we take the exponential of the $\log(\text{odds})$. Based on the parameter estimates in Example 18.4,

$$\text{odds} = e^{-0.5215+0.8431x} = e^{-0.5215} \times e^{0.8431x}$$

From this, the ratio of the odds for a server wearing red ($x = 1$) and for a server wearing a different color ($x = 0$) is

$$\frac{\text{odds}_{\text{red}}}{\text{odds}_{\text{other}}} = e^{0.8431} = 2.324 \blacksquare$$

The transformation $e^{0.8431}$ undoes the natural logarithm and transforms the logistic regression slope into an **odds ratio**, in this case, the comparison of odds that a male customer tips when a server is wearing red to the odds that a male customer tips when a server is wearing a different color. In other words, we can multiply the odds of tipping when a server wears a different color by the odds ratio to obtain the odds of tipping for a server wearing red:

$$\text{odds}_{\text{red}} = 2.324 \times \text{odds}_{\text{other}}$$

In this case, the odds of tipping when a server wears red are about 2.3 times the odds when a server wears a different color.

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 servers wearing a different color as 1 and servers wearing red as 0, the sign of the slope would be reversed, and the fitted equation would be $\log(\text{odds}) = 0.3216 - 0.8431x$, and the odds ratio would be $e^{-0.8431} = 0.430$. The odds of tipping for servers wearing a different color are roughly 43% of the odds for servers wearing red.

Of course, it is often the case that the explanatory variable is quantitative rather than an indicator variable. We must then use software to fit the logistic regression model. Here is an example.

EXAMPLE 18.6

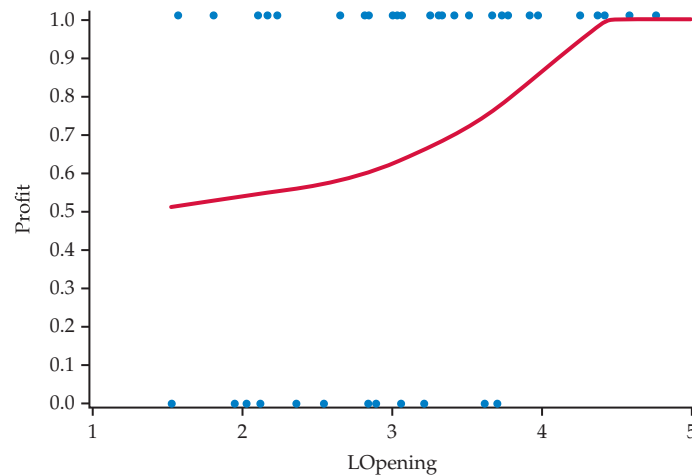


PROFIT

Will a Movie Be Profitable? The MOVIE data set (described on page 635) includes both the movie's budget and the total U.S. revenue. For this example, we classify each movie as "profitable" ($y = 1$) if U.S. revenue is larger than the budget and "unprofitable" ($y = 0$) otherwise. This is our response variable.

The data set contains several explanatory variables, but we focus here on the natural logarithm of the opening-weekend revenue, LOpening . Figure 18.2

FIGURE 18.2 Scatterplot of the movie profit data with a scatterplot smoother, for Example 18.6. The smoother suggests the upper half of an S-shape similar to those shown in Figure 18.1.



is a scatterplot of the data with a scatterplot smoother (page 69). The probability that a movie is profitable increases with the log opening-weekend revenue. Because an S-shaped curve like those in Figure 18.1 is suggested by the smoother, we fit the logistic regression model

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

where p is the probability that the movie is profitable and x is the log opening-weekend revenue. The model for estimated log odds fitted by software is

$$\log(\text{odds}) = b_0 + b_1 x = -1.73 + 0.827x$$

The estimated odds ratio is $e^{b_1} = 2.286$. Because opening-weekend revenue is on the log scale, this means that if opening-weekend revenue x were roughly $e^1 = 2.72$ times larger (e.g., \$18.1 million to \$49.2 million), the odds that the movie will be profitable increase by 2.3 times. ■

APPLY YOUR KNOWLEDGE

18.5 Fitted model for energy drink commercials. Refer to Exercises 18.1 and 18.3. Find the estimates b_0 and b_1 and give the fitted logistic model. What is the estimated odds ratio for a female to choose Commercial A ($x = 1$) versus a male to choose Commercial A ($x = 0$)?

18.6 Fitted model for use of audio/video sharing. Refer to Exercises 18.2 and 18.4. Find the estimates b_0 and b_1 and give the fitted logistic model. What is the odds ratio for small ($x = 1$) versus large ($x = 0$) companies?

The odds ratio interpretation of the estimated slope parameter is a very attractive feature of the logistic regression model. The health sciences, for example, have used this model extensively to identify risk factors for disease and illness. There are other statistical models, such as **probit regression**, that describe binary responses, but none of them has this interpretation.

SECTION 18.1 SUMMARY

- **Logistic regression** explains a success-or-failure response variable in terms of at least one explanatory variable.
- If p is a proportion of successes, then the **odds** of a success are $p/(1-p)$, the ratio of the proportion of successes to the proportion of failures.

- The **simple logistic regression model** relates the proportion of successes in the population to one explanatory variable x through the logarithm of the odds (or **logit**) of a success:

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

That is, each value of x gives a different proportion p of successes. The **data** are n values of x , with observed success or failure for each. The model assumes that these n success-or-failure trials are independent, with probabilities of success given by the logistic regression equation.

- The **parameters** of the simple logistic regression model are β_0 and β_1 . Software fits the data to the model, producing estimates b_0 and b_1 of the parameters β_0 and β_1 .
- The **odds ratio** is the ratio of the odds of a success at $x + 1$ to the odds of a success at x . It is found as e^{β_1} , where β_1 is the slope in the logistic regression equation.

SECTION 18.1 EXERCISES

For Exercises 18.1 and 18.2, see page 18-4; for 18.3 and 18.4, see page 18-5; and for 18.5 and 18.6, see page 18-8.

18.7 What's wrong? For each of the following, explain what is wrong and why.

- If $b_1 = 2$ in a logistic regression analysis, we estimate that the probability of an event is multiplied by 2 when the value of the explanatory variable changes by 1.
- The intercept β_0 is equal to the odds of an event when $x = 0$.
- The odds of an event are 1 minus the probability of the event.

18.8 Interpreting an odds ratio. If we apply the exponential function to the fitted model in Example 18.6, we get

$$\text{odds} = e^{-1.73+0.827x} = e^{-1.73} \times e^{0.827x}$$

Show that for any value of the quantitative explanatory variable x , the odds ratio for increasing x by 1,

$$\frac{\text{odds}_{x+1}}{\text{odds}_x}$$

is $e^{0.827} = 2.286 \approx 2.3$. This justifies the interpretation given at the end of Example 18.6.

18.9 Poor service. In the food service industry, some argue tipping encourages servers to provide discriminate service. If the server expects a good tip, he or she may provide better service. In one survey, 193 servers were surveyed and asked if they ever provided poor service because they did not expect a good tip. Ninety-six replied yes.³

- What proportion of the servers have provided poor service because of an expected bad tip?
- What are the odds that a server will have provided bad service given an expected bad tip?

(c) What proportion of the servers did not provide bad service?

(d) What are the odds that a server will not have provided bad service?

(e) How are your answers to parts (b) and (d) related?

18.10 Is a movie profitable? In Example 18.6 (page 18-7), we developed a model to predict whether a movie will be profitable based on log opening-weekend revenue. What are the predicted odds of a movie being profitable if the opening-weekend revenue is

- \$20 million?
- \$35 million?
- \$50 million?

18.11 Converting odds to probability. Refer to the previous exercise. For each opening-weekend revenue, compute the estimated probability that the movie is profitable.

18.12 How do millennials invest? A PNC Investments survey compared the investing habits of female and male millennials.⁴ Their focus was on millennials with self-reported investments of \$5000 or more or a qualified retirement plan of at least \$1000. One question asked whether they are confident they're saving enough for the future. Here are the data:

Confident	Sex	
	Female	Male
No	410	346
Yes	144	230

- What percent of males and what percent of females are confident they are saving enough?

- (b) Compute the log odds for the males and the log odds for the females feeling confident they're saving enough.
- (c) Write the logistic regression model for this problem using the log odds of feeling confident as the response variable and an indicator for males as the explanatory variable.
- (d) Using the model in (c) and your logits in (b), find the estimates b_0 and b_1 .
- (e) What is the estimated odds ratio for a male-investing millennial to feel confident in saving enough for the future versus a female-investing millennial to feel confident?

18.13 The mobile gender gap. A recent report stated that the mobile ownership gender gap in South Asia is 26%.⁵ Here are the data for India:

Own Mobile	Sex	
	Female	Male
No	343	168
Yes	637	880

- (a) What percent of males and what percent of females own a mobile phone?
- (b) Is the difference in percents close to the reported gap for all of South Asia? Explain your answer.
- (c) What is the estimated odds ratio for a male in India to own a mobile phone versus a female in India to own a mobile phone?
- (d) Use your estimate in (b) to find the estimate of β_1 for a logistic regression using an indicator for males as the explanatory variable.

18.14 Business travel. The Best Western Small Business Travel survey reported that 355 of 400 U.S. small business owners plan as many business trips this fall as last year.⁶

- (a) What proportion of U.S. small business owners plan as many trips as last year?
- (b) What are the odds that an owner will say that his or her company plans as many business trips as last year?
- (c) What proportion of owners said that they do not plan as many trips this year?
- (d) What are the odds that an owner will say that they are cutting back on business trips this year?
- (e) How are your answers to parts (b) and (d) related?

18.15 Know your customers. To devise effective marketing strategies, it is helpful to know the characteristics of your customers. A study compared demographic characteristics of people who use the Internet for travel arrangements and of people who do not.⁷ Of 1132 Internet users, 643 had completed college. Among the 852 nonusers, 349 had completed college. Model the log odds of using the Internet to make travel arrangements with an indicator variable for having completed college as the explanatory variable. Summarize your findings.

18.16 Does income relate to use of the Internet? The study mentioned in the previous exercise also asked about income. Among Internet users, 493 reported income of less than \$50,000, and 378 reported income of \$50,000 or more. (Not everyone answered the income question.) The corresponding numbers for nonusers were 477 and 200. Repeat the analysis using an indicator variable for income of \$50,000 or more as the explanatory variable. What do you conclude?

18.2 Inference for Logistic Regression

When you complete this section, you will be able to:

- Identify the estimates of the simple logistic regression parameters in software output and write the equation for the fitted model.
- Use the estimate and standard error of the regression slope to construct a confidence interval or perform a significance test for the null hypothesis that the slope is zero.
- Compute and interpret the odds ratio and the 95% confidence interval for the odds ratio.

Statistical inference for logistic regression with one explanatory variable is 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.

Wald statistic

The statistic z is sometimes called the **Wald statistic**. Output from some statistical software reports the significance test result in terms of the square of the z statistic.

$$X^2 = z^2$$

chi-square statistic,
p. 549



This statistic is called a chi-square statistic. When the null hypothesis is true, it has a distribution that is approximately a χ^2 distribution with one degree of freedom, and the P -value is calculated as $P(\chi^2 \geq X^2)$. Because the square of a standard Normal random variable has a χ^2 distribution with one degree of freedom, the z statistic and the chi-square statistic give the same results for statistical inference.

CONFIDENCE INTERVALS AND SIGNIFICANCE TESTS FOR LOGISTIC REGRESSION

An approximate **level C confidence interval for the slope** β_1 in the logistic regression model 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** e^{β_1} . A **level C confidence interval for the odds ratio** 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 standard Normal critical value 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 χ^2 having the χ^2 distribution with one degree of freedom, the P -value for a test of H_0 against $H_a: \beta_1 \neq 0$ is approximately $P(\chi^2 \geq X^2)$.

We have expressed the null hypothesis in terms of the slope β_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 18.7



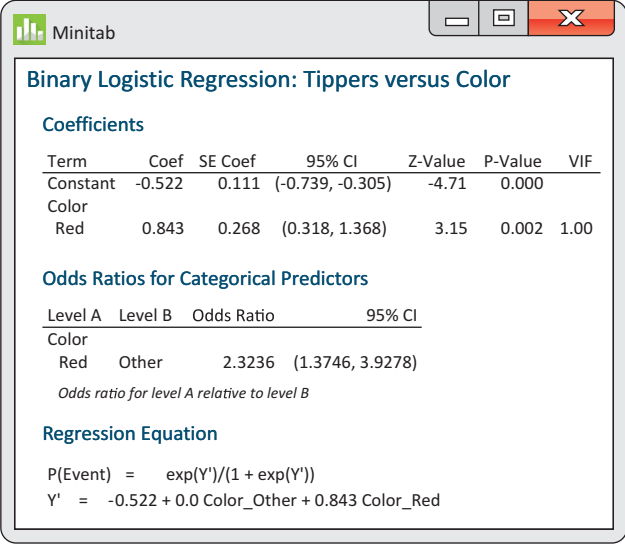
RED



Computer Output for Tipping Study Figure 18.3 gives the output from Minitab and JMP for the tipping study. The parameter estimates match those we calculated in Example 18.4. The standard errors are 0.1107 and 0.2678. A 95% confidence interval for the slope is

$$\begin{aligned} b_1 \pm z^* \text{SE}_{b_1} &= 0.8431 \pm (1.96)(0.2678) \\ &= 0.8431 \pm 0.5249 \end{aligned}$$

FIGURE 18.3 Logistic regression output from Minitab and JMP for the tipping behavior data, for Example 18.7.



Binary Logistic Regression: Tippers versus Color

Coefficients

Term	Coef	SE Coef	95% CI	Z-Value	P-Value	VIF
Constant	-0.522	0.111	(-0.739, -0.305)	-4.71	0.000	
Color						
Red	0.843	0.268	(0.318, 1.368)	3.15	0.002	1.00

Odds Ratios for Categorical Predictors

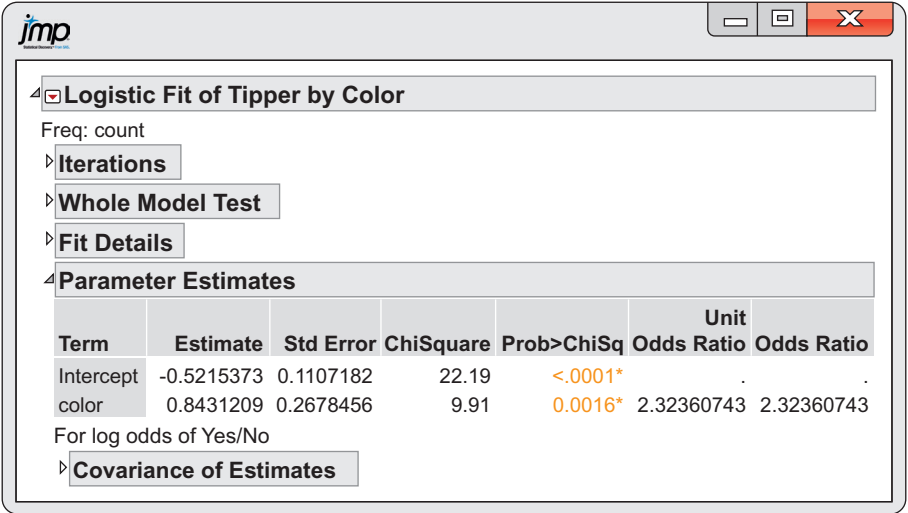
Level A	Level B	Odds Ratio	95% CI
Color			
Red	Other	2.3236	(1.3746, 3.9278)

Odds ratio for level A relative to level B

Regression Equation

$P(\text{Event}) = \frac{\exp(Y')}{1 + \exp(Y')}$

$Y' = -0.522 + 0.0 \text{ Color_Other} + 0.843 \text{ Color_Red}$



Logistic Fit of Tipper by Color

Freq: count

- Iterations
- Whole Model Test
- Fit Details
- Parameter Estimates

Term	Estimate	Std Error	ChiSquare	Prob>ChiSq	Unit	Odds Ratio
Intercept	-0.5215373	0.1107182	22.19	<.0001*		
color	0.8431209	0.2678456	9.91	0.0016*	2.32360743	2.32360743

For log odds of Yes/No

- Covariance of Estimates


We are 95% confident that the slope is between 0.3182 and 1.3680. Both Minitab and JMP provide the odds ratio estimate. Minitab also provides the 95% confidence interval. If this interval is not provided, it is easy to compute from the interval for the slope β_1 :


$$(e^{b_1 - z^* SE_{b_1}}, e^{b_1 + z^* SE_{b_1}}) = (e^{0.3182}, e^{1.3680}) \\ = (1.375, 3.927)$$

We conclude, “Servers wearing red are more likely to be tipped than servers wearing a different color (odds ratio = 2.324, 95% CI = 1.375 to 3.928).” ■

It is standard to use 95% confidence intervals, and software often reports these intervals. A 95% confidence interval for the odds ratio also provides a test of the null hypothesis that the odds ratio is 1 at the 5% significance level. If the confidence interval does not include 1, we reject H_0 and conclude that the odds for the two groups are different; if the interval does include 1, the data do not provide enough evidence to distinguish the groups in this way.

APPLY YOUR KNOWLEDGE

18.17 Inference for energy drink commercials. Use software to run a logistic regression analysis for the energy drink commercial data of Exercise 18.1 (page 18-4). Summarize the results of the inference.  ENERGY

18.18 Inference for audio/visual sharing. Use software to run the logistic regression analysis for the audio/visual sharing data of Exercise 18.2 (page 18-4). Summarize the results of the inference.  AVSHARE



Examples of logistic regression analyses

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

EXAMPLE 18.8



PEST

Effectiveness of an Insecticide As part of a cost-effectiveness study, a wholesale florist company ran an experiment to examine how well the insecticide rotenone kills an aphid called *Macrosiphoniella sanborni* that feeds on the chrysanthemum plant.⁸ The explanatory variable is the concentration (in log of milligrams per liter) of the insecticide. About 50 aphids each were exposed to one of five concentrations. Each insect was either killed or not killed. Here are the data, along with the results of some calculations:

Concentration x (log scale)	Number of insects	Number killed	Proportion killed \hat{p}	Log odds
0.96	50	6	0.1200	-1.9924
1.33	48	16	0.3333	-0.6931
1.63	46	24	0.5217	0.0870
2.04	49	42	0.8571	1.7918
2.32	50	44	0.8800	1.9924

Because there are replications at each concentration, we can calculate the proportion killed and estimate the log odds of death at each concentration. The logistic model in this case assumes that the log odds are *linearly* related to log concentration. Least-squares regression of log odds on log concentration gives the fit illustrated in Figure 18.4. There is a clear linear relationship, which justifies our use of the logistic model. The logistic regression fit for the proportion killed appears in Figure 18.5. It is a transformed version of Figure 18.4 with the fit calculated using the logistic model rather than least squares. ■

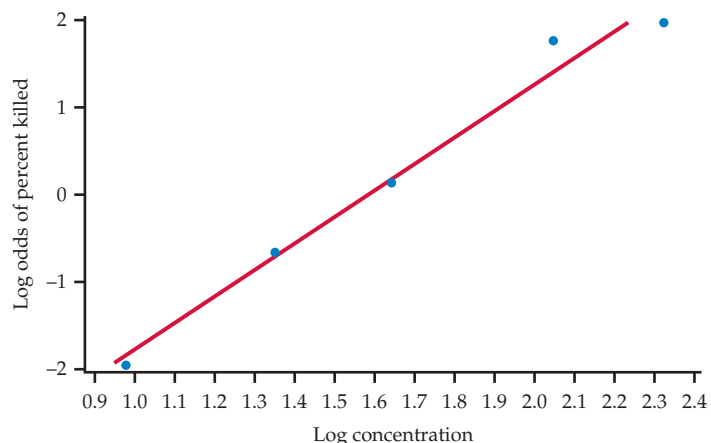
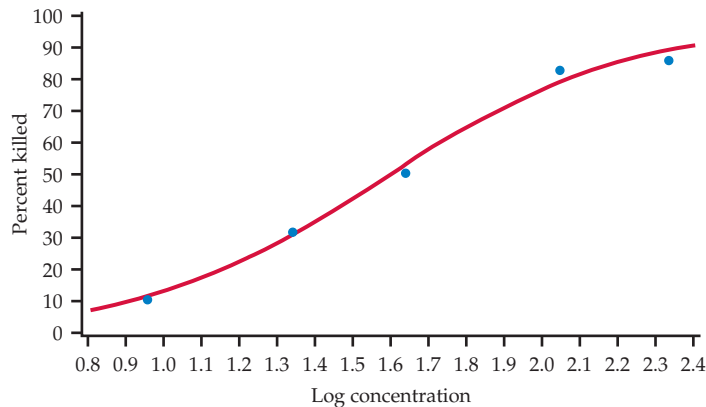


FIGURE 18.4 Plot of log odds of percent killed versus log concentration for the insecticide data, for Example 18.8.

FIGURE 18.5 Plot of the percent killed versus log concentration with the logistic fit for the insecticide data, for Example 18.8.



When the explanatory variable has several values, we can often use graphs like those in Figures 18.4 and 18.5 to visually assess whether the logistic regression model seems appropriate. Just as a scatterplot of y versus x in simple linear regression should show a linear pattern, a plot of log odds versus x in logistic regression should be close to linear. Just as in simple linear regression, outliers in the x direction should be avoided because they may overly influence the fitted model.

The graphs strongly suggest that insecticide concentration affects the kill rate in a way that fits the logistic regression model. Is the effect statistically significant? 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 is not small, our graph may be misleading. As usual, we must add inference to our data analysis.

EXAMPLE 18.9



PEST1

Does Concentration Affect the Kill Rate? Figure 18.6 gives the output from JMP and Minitab for logistic regression analysis of the insecticide data. 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, and 2.32. From the JMP output, we see that the fitted model is

$$\log(\text{odds}) = b_0 + b_1 x = -4.8923 + 3.1088x$$

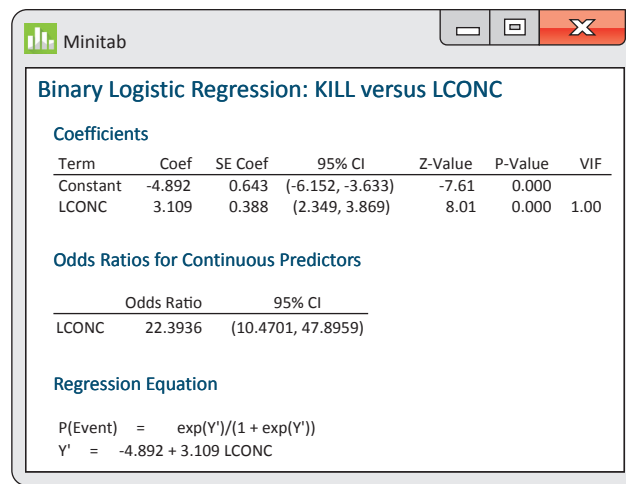
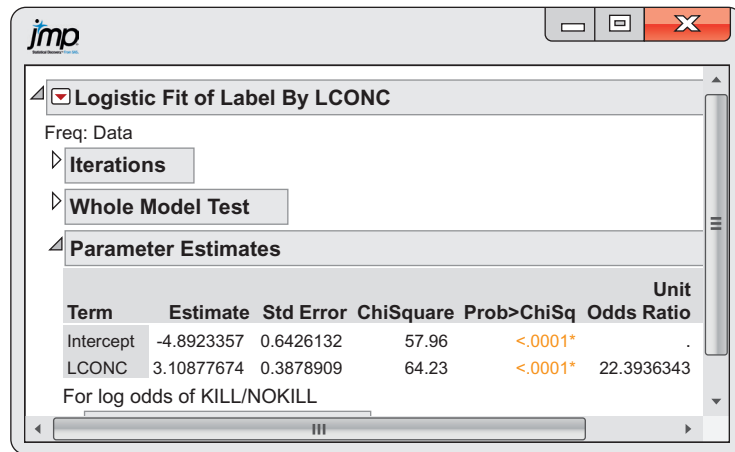
or

$$\frac{\hat{p}}{1-\hat{p}} = e^{-4.8923+3.1088x}$$

Figure 18.5 is a graph of the fitted \hat{p} given by this equation against x , along with the data used to fit the model. JMP gives the statistic X^2 under the heading “ChiSquare.” The null hypothesis that $\beta_1 = 0$ is clearly rejected ($X^2 = 64.23$, $P < 0.0001$).

The estimated odds ratio is 22.394. 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 given in the Minitab output: (10.470, 47.896).

FIGURE 18.6 Logistic regression output from JMP and Minitab for the insecticide data, for Example 18.9.



Remember that the test of the null hypothesis that the slope is 0 is the same as the test of the null hypothesis that the odd ratio is 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 ($X^2 = 64.23$, $P < 0.0001$; 95% CI = 10.5 to 47.9).” ■

APPLY YOUR KNOWLEDGE

18.19 Find the 95% confidence interval for the slope. Using the information in the output of Figure 18.6, find a 95% confidence interval for β_1 .

18.20 Find the 95% confidence interval for the odds ratio. Using the estimate b_1 and its standard error in the output of Figure 18.6, find the 95% confidence interval for the odds ratio and verify that this agrees with the interval given by Minitab.

18.21 X^2 or z . The Minitab output in Figure 18.6 does not give the value of X^2 . The column labeled “Z-Value” provides similar information.

(a) Find the value under the heading “Z-Value” for the predictor LCONC. Verify that this value is simply the estimated coefficient divided by its standard error. This is a z statistic that has approximately the standard Normal distribution if the null hypothesis (slope 0) is true.

(b) Show that the square of z is X^2 . The two-sided P -value for z is the same as P for X^2 .

In Example 18.6, we studied the problem of predicting whether a movie will be profitable using the log opening-weekend revenue as the explanatory variable. We now revisit this example to include the results of inference.

EXAMPLE 18.10



PROFIT

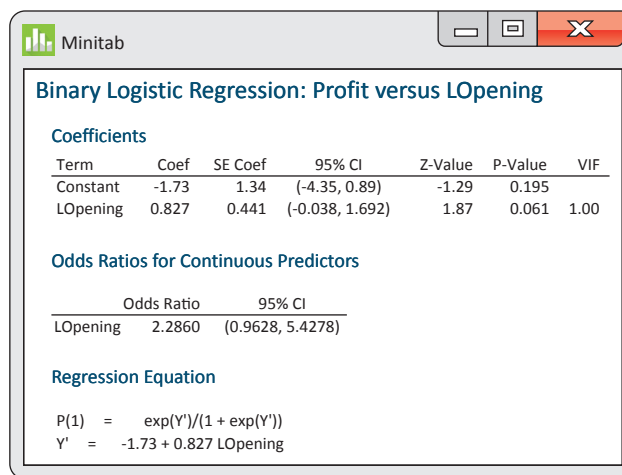
Predicting a Movie's Profitability Figure 18.7 gives the output from Minitab for a logistic regression analysis using log opening-weekend revenue as the explanatory variable. The fitted model is

$$\log(\text{odds}) = b_0 + b_1x = -1.73 + 0.827x$$

This agrees with the result reported in Example 18.6.

From the output, we see that because $P = 0.061$, we cannot reject the null hypothesis that the slope $\beta_1 = 0$ at the 5% significance level. The value of the test statistic is $z = 1.87$, calculated from the estimate $b_1 = 0.827$ and its standard error $SE_{b_1} = 0.441$. Minitab reports the odds ratio as 2.286, with a 95% confidence interval of (0.9628, 5.4278). Notice that this confidence interval contains the value 1, which is another way to assess $H_0: \beta_1 = 0$. In this case, we don't have enough evidence to conclude that this explanatory variable, by itself, is helpful in predicting the probability that a movie will be profitable. ■

FIGURE 18.7 Logistic regression output from Minitab for the movie profit data with log opening-weekend revenue as the explanatory variable, for Example 18.10.



We estimate that a one-unit increase in the log opening-weekend revenue will increase the odds that the movie is profitable about 2.3 times. The data, however, do not give us a very accurate estimate. We do not have strong enough evidence to conclude that movies with higher opening-weekend revenues are more likely to be profitable. Establishing the true relationship accurately would require more data.

SECTION 18.2 SUMMARY

- Software fits the data to the model, producing estimates b_0 and b_1 of the parameters β_0 and β_1 . Software also produces standard errors for these estimates.
- A level C confidence interval for the slope β_1 is

$$b_1 \pm z^* SE_{b_1}$$

A **level C confidence interval for the odds ratio** e^{β_1} is obtained by transforming the confidence interval for the slope,

$$(e^{b_1 - z^* SE_{b_1}}, e^{b_1 + z^* SE_{b_1}})$$

In these expressions, z^* is the standard Normal critical value with area C between $-z^*$ and z^* .

- The null hypothesis that x does not help predict p in the logistic regression model is $H_0: \beta_1 = 0$ or $H_0: e^{\beta_1} = 1$ in terms of the odds ratio. To test this hypothesis, compute the **test statistic**

$$X^2 = \left(\frac{b_1}{SE_{b_1}} \right)^2$$

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

SECTION 18.2 EXERCISES

For Exercises 18.17 and 18.18, see page 18-13; and for 18.19 to 18.21, see page 18-15.

18.22 Using Table F or software. For each of the following situations, use the information provided to compute the P -value for the test of $H_0: \beta_1 = 0$.

- $n = 100$, $b_1 = 0.42$, $X^2 = 5.31$.
- $n = 75$, $b_1 = 1.2$, $SE_{b_1} = 0.43$.
- $n = 200$, $b_1 = -2.3$, $z = -4.2$.

18.23 What's wrong? For each of the following statements, explain what is wrong and why.

- To test the hypothesis $H_0: \beta_1 = 0$ in a simple logistic regression involving n observations, we use χ^2 distribution with $n - 1$ degrees of freedom.
- The confidence interval for the odds ratio is $(e^{b_0 - z^* SE_{b_0}}, e^{b_0 + z^* SE_{b_0}})$.
- We take the square root of the Wald statistic to get the chi-square statistic.

18.24 Following brands through social media.

PricewaterhouseCoopers (PwC) surveyed 1000 online shoppers in the United States and China.⁹ One question asked if the online shopper followed brands they purchased through social media. Here are the results:

Country	Social Media	
	No	Yes
United States	324	676
China	403	597

- What are the proportions of online shoppers who follow brands through social media in each country?

- What is the odds ratio for comparing U.S. online shoppers with Chinese online shoppers?

- Write the logistic regression model for this problem using the log odds of following brands through social media as the response variable and country as an indicator explanatory variable (U.S. = 1).

- Software gives the estimated slope $b_1 = -2.5043$ and its standard error $SE_{b_1} = 0.1377$. Transform this result to the odds scale and compare it with your answer in part (b).

- Construct a 95% confidence interval for the odds ratio and write a short conclusion.

18.25 Analysis of a reduction in force. To meet competition or cope with economic slowdowns, corporations sometimes undertake a “reduction in force” (RIF), in which 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:

Terminated	Over 40	
	No	Yes
Yes	17	71
No	564	835

- Write the logistic regression model for this problem using the log odds of a termination as the response variable and an indicator for over 40 years of age as the explanatory variable.

(b) Explain the assumption 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 $b_1 = 1.0371$ and its standard error $SE_{b_1} = 0.2755$. Transform the results to the odds scale. Summarize the results and write a short conclusion.

18.26 Marketability of tomatoes. A study was conducted to model the marketability of tomato fruit harvested at different maturity stages, disinfection treatments, and storage times.¹⁰ Fruit firmness, measured in newtons (N), was one predictor variable the researchers considered. Each fruit was initially measured for firmness, then stored for 28 days, and finally scored as being marketable or not.

(a) Software gives the estimated intercept $b_0 = -2.243$ with standard error $SE_{b_0} = 0.4573$ and an estimated slope $b_1 = 0.244$ with standard error $SE_{b_1} = 0.1377$. Write the estimated logistic regression model and describe what it tells us about the relationship between the probability of a fruit being marketable and its firmness 28 days earlier.

(b) Test whether fruit firmness helps predict the probability of being marketable.


(c) Construct a 95% confidence interval for the odds ratio and write a short conclusion.


(d) Create a plot with firmness on the x axis and the estimated probability of marketability on the y axis. Use the range of the study firmness measures, specifically 3.8 N to 34.3 N.

(e) A naive tomato farmer decides to pick all his tomatoes so that the firmness is close to 45 N. He argues that this guarantees all his tomatoes will be marketable. Explain what is wrong with his logic.

18.27 Deli refrigerators. The FDA recommends all retail deli refrigerators be kept $<41^\circ\text{C}$. This provision is designed to control bacterial growth on foods such as deli meat. A study looked at the relationship between various food safety practices and deli characteristics and the event that a deli refrigerator exceeds 41°C .¹¹ One predictor of interest was the number of refrigerators at the deli. The following table summarizes the results.

	Number of refrigerators	
	1	> 1
Exceed 41°C		
Yes	16	26
No	118	85

Analyze the data using $\alpha = 0.05$ significance and summarize your conclusions. Make sure to include a 95% confidence interval for the odds ratio.  **FRIDGE**

18.28 Health club membership. A local health club is considering raising its monthly membership fee. To help in this decision, it decides to poll a random sample of $n = 50$ members. Each member was told a new monthly rate, ranging from \$30 to \$60 a month, and asked if they would remain a member if the club changed their rate to this amount. The current rate is \$35 a month. Analyze the data using $\alpha = 0.05$ significance and summarize your conclusions.  **HCLUB**

18.3 Multiple Logistic Regression

When you complete this section, you will be able to:

- Describe the statistical model for logistic regression with multiple explanatory variables.
- Identify the estimates of the regression parameters in software output and write the equation for the fitted model.
- Use software output to test the null hypothesis that all regression slopes are zero.
- Use software output to obtain confidence intervals for the regression coefficients or significance test results for each regression coefficient.
- Interpret the odds ratio and its confidence interval for each explanatory variable.

The PROFIT data set includes several explanatory variables. Example 18.10 examines the model where log opening-weekend revenue alone is used to predict the odds that the movie will have a total U.S. box office revenue greater than the movie budget. Perhaps combining log opening-weekend revenue with other explanatory variables will give us a helpful prediction. We use **multiple logistic regression** to investigate this. 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 13. The statistical concepts are similar, although the computations are now more complex. Here is the analysis.

EXAMPLE 18.11



ANOVA F test,
p. 641

Multiple Logistic Regression As in Example 18.10, we predict the odds that a movie will be profitable. The explanatory variables are log opening-weekend revenue (LOpening), the length of the movie (Minutes), and the movie rating (Rating1). For the movie rating, we use an indicator variable

$$\text{Rating1} = \begin{cases} 1 & \text{if the rating is PG-13 or R} \\ 0 & \text{if the rating is G or PG} \end{cases}$$

Figure 18.8 gives the JMP output. From the output, we see that the fitted model is

$$\begin{aligned} \log(\text{odds}) &= b_0 + b_1 \text{LOpening} + b_2 \text{Minutes} + b_3 \text{Rating1} \\ &= 2.606 + 1.499 \text{LOpening} - 0.060 \text{Minutes} - 0.435 \text{Rating1} \end{aligned}$$

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

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

is tested by a chi-square statistic with three degrees of freedom. JMP reports this to be 8.47 with a *P*-value of 0.0372. We reject H_0 and conclude that one or more of the explanatory variables can be used to predict the odds that the movie is profitable.

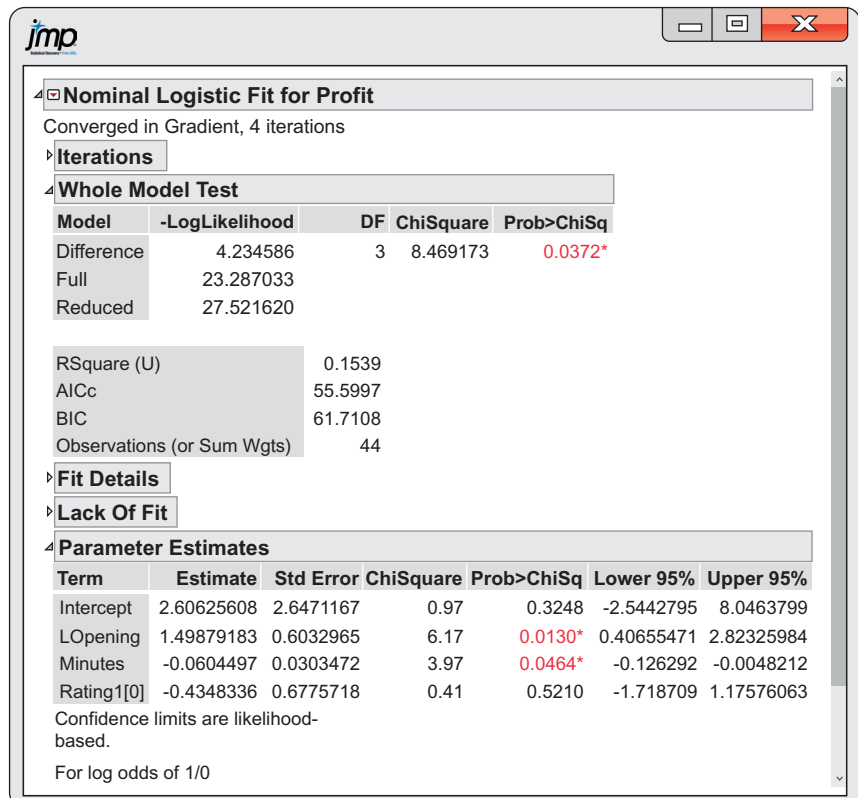


FIGURE 18.8 Multiple logistic regression output from JMP for the movie profit data with log opening-weekend revenue, movie length, and movie rating as the explanatory variables, for Examples 18.11 and 18.12.

Next, we examine the coefficients for each variable and the tests that each of these is 0 *in a model that contains the other two*. The P -values are 0.0130, 0.0464, and 0.5210. The null hypothesis $H_0: \beta_3 = 0$ cannot be rejected. That is, log opening-weekend revenue and the movie's length add significant predictive ability once the other two explanatory variables are already in the model. ■

Because the explanatory variables are correlated, however, we cannot conclude that log opening-weekend revenue and the movie's length make up the best predictive model. 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.

To help with the interpretation of the results, we can transform these estimates to the odds scale, just as we did in simple logistic regression. It is important to note, however, that the resulting ratios and confidence intervals refer to the change in odds for a unit change in a specific explanatory variable, assuming all other explanatory variables do not change. Some call these the **adjusted odds ratios** because the odds ratios are obtained from estimates that account for the other explanatory variables in the model.

EXAMPLE 18.12

Transforming the Estimates to the Odds Scale To get to the odds scale, we take the exponential of each regression coefficient for an explanatory variable. In Figure 18.8, we are also provided the 95% confidence intervals for each regression coefficient. To get the 95% confidence intervals for these odds ratios, we simply take the exponential of each endpoint, just like we did in simple logistic regression.

The regression estimate for log opening-weekend revenue is 1.499. Taking the exponential results in an odds ratio of 4.48. This value is much larger than the one obtained using sample logistic regression in Example 18.10 (page 18-16). We now estimate that a one-unit increase in the log opening-weekend revenue will increase the odds that the movie is profitable about 4.5 times.

The 95% confidence interval is

$$(e^{0.4066}, e^{2.8233}) = (1.50, 16.83)$$

Notice that this confidence interval no longer contains the value 1, which is another way to assess $H_0: \beta_1 = 0$, given the other explanatory variables in the model. ■

SECTION 18.3 SUMMARY

- In **multiple logistic regression** the response variable has two possible values, as in logistic regression, but there is more than one explanatory variable.
- As in multiple regression, there is an **overall test** for all the explanatory variables. The null hypothesis that the coefficients for all the explanatory variables are zero is tested by a statistic that has a distribution that is approximately χ^2 with degrees of freedom equal to the number of explanatory variables. The P -value is approximately $P(\chi^2 \geq X^2)$.
- Hypotheses about **individual coefficients**, $H_0: \beta_j = 0$ or $H_0: e^{\beta_j} = 1$ in terms of the odds ratio, are tested by a statistic that is approximately χ^2 with 1 degree of freedom. The P -value is approximately $P(\chi^2 \geq X^2)$. As in multiple regression, these tests assess the contribution of each explanatory variable given the other explanatory variables are already in the model.

SECTION 18.3 EXERCISES

18.29 Tipping behavior in Canada. The Consumer Report on Eating Share Trends (CREST) contains data that cover all provinces of Canada and that describe away-from-home food purchases by roughly 4000 households per quarter. Researchers recently restricted their attention to restaurants at which tips would normally be given.¹² From a total of 73,822 observations, high and low tipping variables were created based on whether the observed tip rate was above 20% or below 10%, respectively. They then used logistic regression to identify explanatory variables associated with either high or low tips. Here is a table summarizing what they termed the stereotype-related variables for the high-tip analysis:

Explanatory variable	Odds ratio
Senior adult	0.7420*
Sunday	0.9970
English as second language	0.7360*
French-speaking Canadian	0.7840*
Alcoholic drinks	1.1250*
Lone male	1.0220

The starred odds ratios were significant at the 0.01 level. Write a short summary explaining these results in terms of the odds of leaving a high tip.

18.30 Finding the best model. In Example 18.11 (page 18-19), we looked at a multiple logistic regression for movie profitability based on three explanatory variables. Complete the analysis by looking at the three models with two explanatory variable models and the three models with single variables. Create a table that includes the parameter estimates and their P -values as well as the X^2 statistic and degrees of freedom. Based on the results, which model do you think is the best? Explain your answer.



18.31 E-government use in Canada. Electronic government (e-government) provides digital means, such as an email address or a website, for citizens to contact public officials. The vision behind e-government is to create a more citizen-focused government. One study used survey data to determine what factors are related to a citizen using an e-government website rather than visiting or calling a government office.¹³ The dependent variable refers to whether the citizen used the website. Explanatory variables include sex (1 = female, 0 = male), daily Internet use (1 = yes, 0 = no), age (six ordered categories numbered 1 through 6), household income (seven ordered categories numbered 1 through 7), size of the community (six ordered categories numbered 1 through 6), and education (1 = at least some postsecondary education, 0 = other). The following table summarizes the results.

Explanatory variable	Odds ratio
Sex	0.87
Daily Internet use	4.16
Age	0.81
Income	1.01
Community size	0.85
Education	0.97
Intercept	0.66

All but “Education” and “Income” were significant at the 0.05 level.


- Interpret each of the odds ratios in terms of the probability that the individual uses the website.
- Compute the regression coefficients for each of the variables in the table.
- What are the odds that a male college graduate who uses the Internet daily, and is age category 3, household income level 4, and community size 5 is using the Internet?

18.32 CEO overconfidence/dominance and corporate acquisitions. The acquisition literature suggests that takeovers occur either due to conflicts between managers and shareholders or to create a new entity that exceeds the sum of its previously separate components. Other research has offered managerial hubris as a third option, but it has not been studied empirically. Recently, some researchers revisited acquisitions over a 10-year period in the Australian financial system.¹⁴ A measure of CEO overconfidence was based on the CEO’s level of media exposure, and a measure of dominance was based on the CEO’s remuneration relative to the firm’s total assets. They then used logistic regression to see whether CEO overconfidence and dominance were positively related to the probability of at least one acquisition in a year. To help isolate the effects of CEO hubris, the model included explanatory variables of firm characteristics and other potentially important factors in the decision to acquire. The following table summarizes the results for the two key explanatory variables:

Explanatory variable	b	$SE(b)$
Overconfidence	0.0878	0.0402
Dominance	1.5067	0.0057

- State the null and alternative hypotheses for each of the explanatory variables.
- Perform the significance tests and determine whether the variables are significant at the 0.05 level.
- Estimate the odds ratio for each variable and construct a 95% confidence interval.
- Write a short summary explaining the results.

18.33 Finding the best model, continued. In Example 18.11 (page 18-19), we looked at a multiple logistic regression for movie profitability based on three

explanatory variables. Now consider all the explanatory variables in the data set to find the best model. Write a short summary describing the model you select.  **PROFIT**

CHAPTER 18 REVIEW EXERCISES

18.34 What's wrong? For each of the following, explain what is wrong and why.

- (a) For a multiple logistic regression with four explanatory variables, the null hypothesis that the regression coefficients of all the explanatory variables are zero is tested with an F test.
- (b) For a logistic regression we assume that the error term in our model has a Normal distribution.
- (c) In logistic regression with two explanatory variables we use a chi-square statistic to test the null hypothesis $H_0: b_1 = 0$ versus a two-sided alternative.

18.35 Stock options. Different kinds of companies compensate their key employees in different ways. Established companies may pay higher salaries, while new companies may offer stock options that will be valuable if the company succeeds. Do high-tech companies tend to offer stock options more often than other companies? One study looked at a random sample of 200 companies. Of these, 91 were listed in the *Directory of Public High Technology Corporations*, and 109 were not listed. Treat these two groups as SRSs of high-tech and non-high-tech companies. Seventy-three of the high-tech companies and 75 of the non-high-tech companies offered incentive stock options to key employees.¹⁵

- (a) What proportion of the high-tech companies offer stock options to their key employees? What are the odds?
- (b) What proportion of the non-high-tech companies offer stock options to their key employees? What are the odds?
- (c) Find the odds ratio using the odds for the high-tech companies in the numerator. Interpret the result in a few sentences.

18.36 Log odds for high-tech and non-high-tech firms. Refer to the previous exercise.

- (a) Find the log odds for the high-tech firms. Do the same for the non-high-tech firms.
- (b) Define an explanatory variable x to have the value 1 for high-tech firms and 0 for non-high-tech firms. For the logistic model, we set the log odds equal to $\beta_0 + \beta_1 x$. Find the estimates b_0 and b_1 for the parameters β_0 and β_1 .
- (c) Show that the odds ratio is equal to e^{b_1} .

18.37 Do the inference. Refer to the previous exercise. Software gives 0.3347 for the standard error of b_1 .

- (a) Find the 95% confidence interval for β_1 .
- (b) Transform your interval in (a) to a 95% confidence interval for the odds ratio.
- (c) What do you conclude?

18.38 Sexual imagery in magazine ads. In what ways do advertisers in magazines use sexual imagery to appeal to youth? One study classified each of 1509 full-page or larger ads as “not sexual” or “sexual,” according to the amount and style of the dress of the male or female model in the ad. The ads were also classified according to the target readership of the magazine.¹⁶ A logistic regression was used to describe the probability that the clothing in the ad was not sexual as a function of several explanatory variables. Here are some of the reported results:

Explanatory variable	b	z test
Reader age	0.50	13.64
Model sex	1.31	72.15
Men's magazines	-0.05	0.06
Women's magazines	0.45	6.44
Constant	-2.32	135.92

Reader age is coded as 0 for young adult and 1 for mature adult. Therefore, the coefficient of 0.50 for this explanatory variable suggests that the probability that the model clothing is *not* sexual is higher when the target reader age is mature adult. In other words, the model clothing is more likely to be sexual when the target reader age is young adult. Model sex is coded as 0 for female and 1 for male. The explanatory variable men's magazines is 1 if the intended readership is men and 0 for women's magazines and magazines intended for both men and women (general interest). The variable women's magazines is coded similarly.

- (a) State the null and alternative hypotheses for each of the explanatory variables.
- (b) Perform the significance tests associated with the z statistics.
- (c) Interpret the sign of each of the statistically significant coefficients in terms of the probability that the model clothing is sexual.
- (d) Write an equation for the fitted logistic regression model.

18.39 Interpret the results. Refer to the previous exercise. The researchers also reported odds ratios with

95% confidence intervals for this logistic regression model. Here is a summary:

Explanatory variable	Odds ratio	95% Confidence limits	
		Lower	Upper
Reader age	1.65	1.27	2.16
Model sex	3.70	2.74	5.01
Men's magazines	0.96	0.67	1.37
Women's magazines	1.57	1.11	2.23

- (a) Explain the relationship between the confidence intervals reported here and the results of the z significance tests that you found in the previous exercise.
- (b) Interpret the results in terms of the odds ratios.
- (c) Write a short summary explaining the results. Include comments regarding the usefulness of the fitted coefficients versus the odds ratios in making a summary.

18.40 Suppose you had twice as many data. Refer to Exercises 18.35 through 18.37. Repeat the calculations assuming that you have twice as many observations with the same proportions. In other words, assume that there are 182 high-tech firms and 218 non-high-tech firms. The numbers of firms offering stock options are 146 for the high-tech group and 150 for the non-high-tech group. The standard error of b_1 for this scenario is 0.2366. Summarize your results, paying particular attention to what remains the same and what is different from what you found in Exercises 18.25 to 18.27.


18.41 How do millennials invest, continued. Refer to Exercise 18.12 (page 18-9). Another question asked whether they feel in complete control of their financial well-being. Here are their responses:

In control	Sex	
	Female	Male
No	377	328
Yes	177	248

- (a) What percent of males and what percent of females feel in complete control of their financial well-being?
- (b) Compute the log odds for the males and the log odds for the females feeling in complete control of their financial well-being.
- (c) Write the logistic regression model for this problem using the log odds of feeling in complete control as the response variable and an indicator for males as the explanatory variable.
- (d) Using the model in (c) and your logits in (b), find the estimates b_0 and b_1 .

- (e) What is the estimated odds ratio for a male to feel in complete control versus a female to feel in complete control?


18.42 Active retail companies versus failed companies.

Case 8.2 (page 429) compared the cash flow of 74 active retail firms with the cash flow of 27 firms that failed. Here we analyze the same data with a logistic regression. The outcome is whether the firm is active, and the explanatory variable is the cash flow. Here is the output from Minitab:  CMPS

Coefficients			
Term	Coef	SE Coef	VIF
Constant	1.068	0.241	
CashFlow	0.0320	0.0114	1.00
Odds Ratios for Continuous Predictors			
	Odds Ratio	95% CI	
CashFlow	1.0326	(1.0098, 1.0558)	

- (a) Give the fitted equation for the log odds that a firm will be active.
- (b) Describe the results of the significance test for the coefficient of cash flow.
- (c) The odds ratio is the estimated amount that the odds of being active would increase when the cash flow is increased by one unit. Report this odds ratio with the 95% confidence interval.
- (d) Write a short summary of this analysis and compare it with the analysis of these data that we performed in Chapter 8. Which approach do you prefer?


For the following five exercises, you will need to construct indicator variables to use categorical variables as explanatory variables in logistic regression. Be sure to review the material in Chapter 13 on models with categorical explanatory variables (pages 655–658) before attempting these exercises.

18.43 Reduction in force using logistic regression. In Exercise 18.25 (page 18-17), hypothetical data are given for a reduction in force (RIF). If there is a statistically significant difference in the RIF proportions based on age group, the employer needs to justify the difference based on other (nondiscriminatory) variables. 


- (a) Run the logistic analysis to predict the odds of being riffed using age group (over 40 years of age or not) as the explanatory variable. Summarize your results.
- (b) What other variables would you add to the model in an attempt to explain the results that you described in part (a)? If these other variables can be shown to be characteristics that relate to job performance, and the age effect is no longer significant in a model that includes these variables, then the analysis provides statistical evidence that can be used to refute a claim of discrimination.

18.44 Sexual imagery in ads. Refer to Exercise 18.38 (page 18-22) concerning the use of sexual imagery in magazine ads. Here is the two-way table of counts for the 1509 ads.

Model dress	Magazine readership			Total
	Women	Men	General interest	
Not sexual	351	514	248	1113
Sexual	225	105	66	396
Total	576	619	314	1509

Use the model dress, expressed as the odds that the dress is sexual, as the response variable and the magazine readership as the explanatory variable. Because there are three magazine readership categories, you will need two indicator variables for this multiple logistic regression analysis. Use the last category, general interest, for the “other” designation when creating these indicator variables.  **IMAGERY**

- A friend has suggested that the three magazine categories be coded as 1, 2, 3 and that this single variable be used as the explanatory variable in the logistic regression. Explain why this analytical strategy is wrong.
- Summarize the results of the significance testing. Do the data support the idea that the sexual content expressed in the model dress varies by the magazine readership?
- Use the estimates for your model and the coding that you used for the explanatory variables to give the estimated log odds for each type of magazine readership.

18.45 Rerun the analysis with a different coding. In the previous exercise, you used the last category, general interest, for the “other” designation when you constructed the indicator variables. Now use the women’s magazine readership as the “other” category and reanalyze the data. Verify that the significance testing results for the effect of the two explanatory variables is the same as in the previous exercise.  **IMAGERY**

18.46 Student athletes and gambling. A survey of student athletes that asked questions about gambling behavior classified students according to the National Collegiate Athletic Association (NCAA) division.¹⁷ For male student athletes, the percent who reported wagering on collegiate sports is given here along with the numbers of respondents in each division:

	Division		
	I	II	III
Percent	17.2	21.0	24.4
Number	5619	2957	4089

- Using the numbers and percents given, calculate the numbers of students who gamble and those who do not for each NCAA division.
- Use two indicator variables to code the explanatory variable, NCAA division. Let the first one be 1 for Division II and 0 otherwise; let the second be 1 for Division III and 0 otherwise. With this coding, the logistic regression model will use the intercept for Division I, the intercept plus the coefficient of the first indicator variable for Division II, and the intercept plus the coefficient of the second indicator variable for Division III.

(c) Run the multiple logistic regression and summarize the results.

18.47 Is there a trend? Refer to the previous exercise. The coding of the indicator variables suggests a way to code models when you expect a pattern in the response that is based on some kind of ordering of the explanatory variable. In some settings this is called detecting a dose response.

- Use the model to give the estimated log odds for each NCAA division.
- Plot these estimates versus division and summarize the results. Does there appear to be a pattern in the results?
- How would you model the pattern that you described in part (b)?

Answers to Odd-Numbered Exercises

- 18.1** For men: the percent who chose Commercial A is $97/220 = 0.4409$, Commercial B is 0.5591. The odds are 0.7886. For women: the percent who chose Commercial A is $90/150 = 0.6$, Commercial B is 0.4. The odds are 1.5.
- 18.3** For men: $\log(0.7886) = -0.2375$. For women: $\log(1.5) = 0.4055$.
- 18.5** $b_0 = -0.2375$; $b_1 = 0.4055 - (-0.2375) = 0.643$. $\log(\text{odds}) = -0.2375 + 0.643x$. The odds ratio is 1.902.
- 18.7** (a) It is not multiplied by 2. When the explanatory variable changes by 1, the odds ratio is increased by a factor of e^2 or 7.389 times. (b) It is missing the log; the intercept is equal to the log odds of an event when $x = 0$. (c) The odds of an event is the probability of the event divided by 1 minus the probability of the event.
- 18.9** (a) $96/193 = 0.4974$. (b) Odds = 0.9897. (c) $97/193 = 0.5026$. (d) Odds = 1.0104. (e) They are reciprocals, (d) = $1/(b)$.
- 18.11** (a) $\frac{2.416}{1 + 2.416} = 0.7073$. (b) $\frac{1.3703}{1 + 1.3703} = 0.7974$. (c) $\frac{5.3724}{1 + 5.3724} = 0.8431$.
- 18.13** (a) For males: 84%; for females: 65%. (b) The difference in percents, which is 1.039, is different from the reported ownership gender gap for South Asia, which is 26%. (c) 2.827. (d) 1.039.
- 18.15** $\log(\text{odds}) = -0.0282 + 0.6393x$. $X^2 = 48.34$; $P\text{-value} < 0.0001$. The odds ratio estimate is 1.8952; those that have completed college are 1.8952 times more likely to use the Internet for travel arrangements than those that have not completed college.
- 18.17** $\log(\text{odds}) = -0.2375 + 0.643x$; the odds ratio is 1.902; the 95% confidence interval is (0.222, 1.064).
- 18.19** (2.349, 3.869).
- 18.21** (a) $Z = 3.109/0.388 = 8.01$. (b) $8.01^2 = 64.16$, which agrees with the output up to rounding error.
- 18.23** (a) In a simple logistic regression, there is 1 degree of freedom for the χ^2 distribution. (b) The confidence interval for the odds ratio is $(e^{b_1 - z^* SE_{b_1}}, e^{b_1 + z^* SE_{b_1}})$. (c) The square root of the Wald statistic is the Z statistic.
- 18.25** (a) $\log(\text{odds}) = -3.5018 + 1.0371x$. (b) The binomial distribution assumes that each employee's termination is independent from another's and the probability of being terminated is the same for each employee. Certainly the latter is not true, as an individual's performance is likely different and largely determines whether he or she is terminated. (c) Odds = $e^{1.0371} = 2.82$, with 95% the confidence interval is (1.644, 4.840). Since the interval does not contain 1, the results are significant at the 5% level. Employees over 40 are 2.82 times more likely to be terminated than those under 40. We could use the additional variables in the logistic regression model to account for their effects before assessing if age has an effect.
- 18.27** $\log(\text{odds}) = -1.998 + 0.8135x$; the 95% confidence interval for the odds ratio is (1.14, 4.46); since the interval does not include 1, the results are significant.
- 18.29** Those that order alcoholic drinks are 12.5% more likely (or 1.125 times as likely) to leave a high tip than those that don't order alcohol. Senior adults are about 25.8% less likely (or 0.742 times as likely) to leave a high tip than those that aren't senior. Those that speak English as a second language are about 26.4% less likely (or 0.736 times as likely) to leave a high tip than their counterparts. Those that are French-speaking Canadians are about 21.6% less likely (or 0.784 times as likely) to leave a high tip than those that aren't French-speaking Canadians.
- 18.31** (a) Females are 0.87 times as likely (13% less likely) to use the website as males. Daily Internet users are 4.16 times as likely to use the website as their counterparts. Older-aged people are less likely to use the website than younger-aged people. Those from larger communities are less likely to use the website than those from smaller communities. Those with different incomes and/or educational attainment are about equally likely to use the website, since they aren't significantly different from 1. (b) Sex: -0.1393; Daily Internet use: 1.4255; Age: -0.2107; Income: 0.01; Size: -0.1625; Education: -0.0305; Intercept: -0.4155. (c) 0.6537.
- 18.33** Answers will vary.
- 18.35** (a) $73/91 = 0.8022$; odds = 4.0556. (b) $75/109 = 0.6881$; odds = 2.2059. (c) Odds ratio = $4.0556/2.2059 = 1.8385$. The high-tech companies are 1.8385 times more likely to offer incentive stock options to key employees than the non-high-tech companies are.

- 18.37** (a) $0.6089 \pm 1.96(0.3347) = (-0.047, 1.265)$.
 (b) (0.954, 3.543). (c) Since the interval in part (b) includes 1, there is no significant difference in the proportions of high-tech and non-high-tech companies that offer stock options to key employees.
- 18.39** (a) If the confidence interval for the odds ratio includes the value 1, the variable is not significant in a logistic regression. (b) Since the Reader age, Model sex, and Women's magazine intervals all do not contain 1, they are all significant. The Men's magazine interval contains 1 and is not significant. (c) Interpreting only significant effects: When the reader age is mature adults, the model clothing is 1.27 to 2.16 times more likely to be not sexual. When the model sex is male, the model clothing is 2.74 to 5.01 times more likely to be not sexual. When the intended readership is women, the model clothing is 1.11 to 2.23 times more likely to be not sexual. The odds ratios are often much easier to interpret than the fitted coefficients.
- 18.41** (a) For males: 43.06%; for females: 31.95%.
 (b) For males: -0.2796 ; for females: -0.7561 .

(c) $\log(\text{odds}) = b_0 + b_1 \text{Males}$. (d) $\log(\text{odds}) = -0.7561 + 0.4765 \text{Males}$.

- 18.43** (a) $\log(\text{odds}) = -3.5017 + 1.0369x$. $X^2 = 14.17$; $P\text{-value} = 0.0002$. The model is significant. For a person over 40, the odds of being RIFFED are $e^{-3.5017+1.0369(1)} = 0.085$. For a person under 40, the odds of being RIFFED are: $e^{-3.5017+1.0369(0)} = 0.030$. (b) Answers will vary. (c) Since the relationship is quite linear, we could use a regression analysis.
- 18.45** The estimated model becomes $\log(\text{odds}) = -0.4447 - 1.1436x_{\text{men}} - 0.8791x_{\text{general}}$. Now both men ($X^2 = 69.7000$; $P\text{-value} < 0.0001$) and general ($X^2 = 29.1872$; $P\text{-value} < 0.0001$) terms are significant.
- 18.47** (a) For Division I: $\log(\text{odds}) = -1.5720 + 0.2472(0) + 0.4416(0) = -1.5720$. For Division II: $\log(\text{odds}) = -1.5720 + 0.2472(1) + 0.4416(0) = -1.3248$. For Division III: $\log(\text{odds}) = -1.5720 + 0.2472(0) + 0.4416(1) = -1.1304$. (b) The plot shows that $\log(\text{odds})$ of gambling increases as Division increases. (c) Since the relationship is quite linear, we could use a regression analysis.