



CHAPTER

10

Correlation and Regression

Objectives

After completing this chapter, you should be able to

- 1 Draw a scatter plot for a set of ordered pairs.
- 2 Compute the correlation coefficient.
- 3 Test the hypothesis $H_0: \rho = 0$.
- 4 Compute the equation of the regression line.
- 5 Compute the coefficient of determination.
- 6 Compute the standard error of the estimate.
- 7 Find a prediction interval.
- 8 Be familiar with the concept of multiple regression.

Outline

Introduction

10-1 Scatter Plots and Correlation

10-2 Regression

10-3 Coefficient of Determination and Standard Error of the Estimate

10-4 Multiple Regression (Optional)

Summary



Statistics Today

Do Dust Storms Affect Respiratory Health?

Southeast Washington state has a long history of seasonal dust storms. Several researchers decided to see what effect, if any, these storms had on the respiratory health of the people living in the area. They undertook (among other things) to see if there was a relationship between the amount of dust and sand particles in the air when the storms occur and the number of hospital emergency room visits for respiratory disorders at three community hospitals in southeast Washington. Using methods of correlation and regression, which are explained in this chapter, they were able to determine the effect of these dust storms on local residents. See Statistics Today—Revisited at the end of the chapter.

Source: B. Hefflin, B. Jalaludin, N. Cobb, C. Johnson, L. Jecha, and R. Etzel, “Surveillance for Dust Storms and Respiratory Diseases in Washington State,” *Archives of Environmental Health* 49, no. 3 (May–June), pp. 170–74. Reprinted with permission of the Helen Dwight Reid Education Foundation. Published by Heldref Publications, 1319 18th St. N.W., Washington, D.C. 20036-1802.

Introduction

In Chapters 7 and 8, two areas of inferential statistics—confidence intervals and hypothesis testing—were explained. Another area of inferential statistics involves determining whether a relationship exists between two or more numerical or quantitative variables. For example, a businessperson may want to know whether the volume of sales for a given month is related to the amount of advertising the firm does that month. Educators are interested in determining whether the number of hours a student studies is related to the student’s score on a particular exam. Medical researchers are interested in questions such as, Is caffeine related to heart damage? or Is there a relationship between a person’s age and his or her blood pressure? A zoologist may want to know whether the birth weight of a certain animal is related to its life span. These are only a few of the many questions that can be answered by using the techniques of correlation and regression analysis. **Correlation** is a statistical method used to determine whether a linear relationship between variables exists. **Regression** is a statistical method used to describe the nature of the relationship between variables, that is, positive or negative, linear or nonlinear.

The purpose of this chapter is to answer these questions statistically:

1. Are two or more variables linearly related?
2. If so, what is the strength of the relationship?

Unusual Stat

A person walks on average 100,000 miles in his or her lifetime. This is about 3.4 miles per day.

3. What type of relationship exists?
4. What kind of predictions can be made from the relationship?

To answer the first two questions, statisticians use a numerical measure to determine whether two or more variables are linearly related and to determine the strength of the relationship between or among the variables. This measure is called a *correlation coefficient*. For example, there are many variables that contribute to heart disease, among them lack of exercise, smoking, heredity, age, stress, and diet. Of these variables, some are more important than others; therefore, a physician who wants to help a patient must know which factors are most important.

To answer the third question, you must ascertain what type of relationship exists. There are two types of relationships: *simple* and *multiple*. In a **simple relationship**, there are two variables—an **independent variable**, also called an explanatory variable or a predictor variable, and a **dependent variable**, also called a response variable. A simple relationship analysis is called *simple regression*, and there is one independent variable that is used to predict the dependent variable. For example, a manager may wish to see whether the number of years the salespeople have been working for the company has anything to do with the amount of sales they make. This type of study involves a simple relationship, since there are only two variables—years of experience and amount of sales.

In a **multiple relationship**, called *multiple regression*, two or more independent variables are used to predict one dependent variable. For example, an educator may wish to investigate the relationship between a student's success in college and factors such as the number of hours devoted to studying, the student's GPA, and the student's high school background. This type of study involves several variables.

Simple relationships can also be positive or negative. A **positive relationship** exists when both variables increase or decrease at the same time. For instance, a person's height and weight are related; and the relationship is positive, since the taller a person is, generally, the more the person weighs. In a **negative relationship**, as one variable increases, the other variable decreases, and vice versa. For example, if you measure the strength of people over 60 years of age, you will find that as age increases, strength generally decreases. The word *generally* is used here because there are exceptions.

Finally, the fourth question asks what type of predictions can be made. Predictions are made in all areas and daily. Examples include weather forecasting, stock market analyses, sales predictions, crop predictions, gasoline price predictions, and sports predictions. Some predictions are more accurate than others, due to the strength of the relationship. That is, the stronger the relationship is between variables, the more accurate the prediction is.

10–1

Scatter Plots and Correlation**Objective 1**

Draw a scatter plot for a set of ordered pairs.

In simple correlation and regression studies, the researcher collects data on two numerical or quantitative variables to see whether a relationship exists between the variables. For example, if a researcher wishes to see whether there is a relationship between number of hours of study and test scores on an exam, she must select a random sample of students, determine the hours each studied, and obtain their grades on the exam. A table can be made for the data, as shown here.

Student	Hours of study x	Grade y (%)
A	6	82
B	2	63
C	1	57
D	5	88
E	2	68
F	3	75

As stated previously, the two variables for this study are called the independent variable and the dependent variable. The independent variable is the variable in regression that can be controlled or manipulated. In this case, the number of hours of study is the independent variable and is designated as the x variable. The dependent variable is the variable in regression that cannot be controlled or manipulated. The grade the student received on the exam is the dependent variable, designated as the y variable. The reason for this distinction between the variables is that you assume that the grade the student earns *depends* on the number of hours the student studied. Also, you assume that, to some extent, the student can regulate or *control* the number of hours he or she studies for the exam.

The determination of the x and y variables is not always clear-cut and is sometimes an arbitrary decision. For example, if a researcher studies the effects of age on a person's blood pressure, the researcher can generally assume that age affects blood pressure. Hence, the variable *age* can be called the *independent variable*, and the variable *blood pressure* can be called the *dependent variable*. On the other hand, if a researcher is studying the attitudes of husbands on a certain issue and the attitudes of their wives on the same issue, it is difficult to say which variable is the independent variable and which is the dependent variable. In this study, the researcher can arbitrarily designate the variables as independent and dependent.

The independent and dependent variables can be plotted on a graph called a *scatter plot*. The independent variable x is plotted on the horizontal axis, and the dependent variable y is plotted on the vertical axis.

A **scatter plot** is a graph of the ordered pairs (x, y) of numbers consisting of the independent variable x and the dependent variable y .

The scatter plot is a visual way to describe the nature of the relationship between the independent and dependent variables. The scales of the variables can be different, and the coordinates of the axes are determined by the smallest and largest data values of the variables.

The procedure for drawing a scatter plot is shown in Examples 10–1 through 10–3.

Example 10–1

Car Rental Companies



Construct a scatter plot for the data shown for car rental companies in the United States for a recent year.

Company	Cars (in ten thousands)	Revenue (in billions)
A	63.0	\$7.0
B	29.0	3.9
C	20.8	2.1
D	19.1	2.8
E	13.4	1.4
F	8.5	1.5

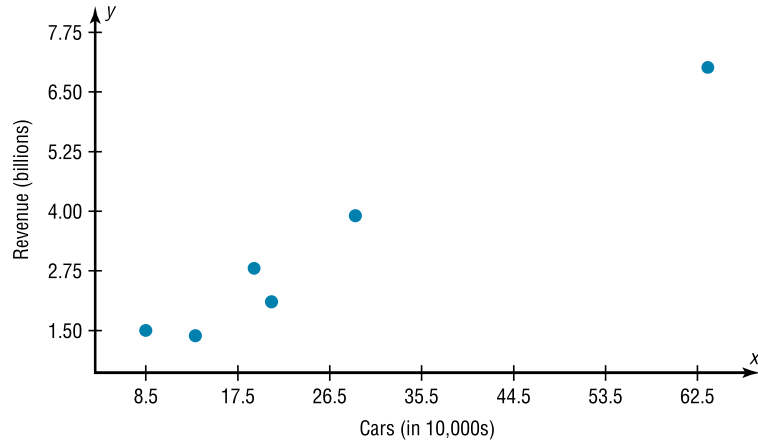
Source: *Auto Rental News*.

Solution

Step 1 Draw and label the x and y axes.

Step 2 Plot each point on the graph, as shown in Figure 10–1.

Figure 10–1
Scatter Plot for
Example 10–1



Example 10–2

Absences and Final Grades



Construct a scatter plot for the data obtained in a study on the number of absences and the final grades of seven randomly selected students from a statistics class. The data are shown here.

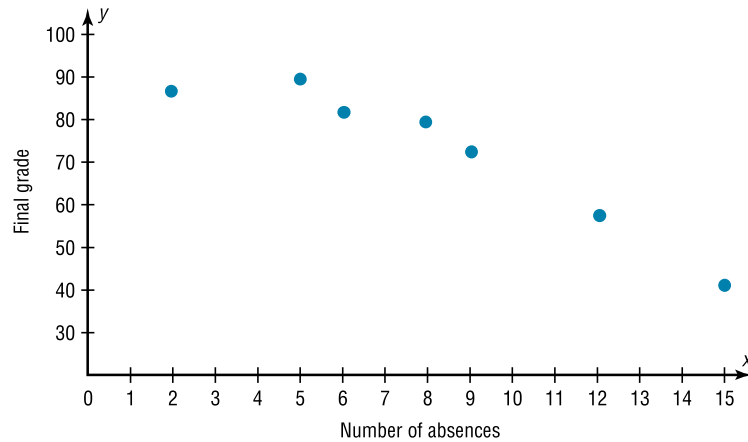
Student	Number of absences x	Final grade y (%)
A	6	82
B	2	86
C	15	43
D	9	74
E	12	58
F	5	90
G	8	78

Solution

Step 1 Draw and label the x and y axes.

Step 2 Plot each point on the graph, as shown in Figure 10–2.

Figure 10–2
Scatter Plot for
Example 10–2



Example 10–3**Age and Wealth**

A researcher wishes to see if there is a relationship between the ages and net worth of the wealthiest people in America. The data for a specific year are shown.

Person	Age x	Net wealth y (\$ billions)
A	73	16
B	65	26
C	53	50
D	54	21.5
E	79	40
F	69	16
G	61	19.6
H	65	19

Source: *Forbes* magazine.

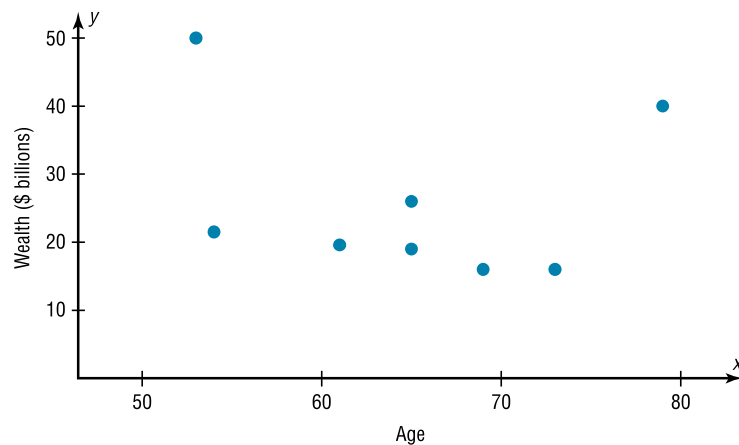
Solution

Step 1 Draw and label the x and y axes.

Step 2 Plot each point on the graph, as shown in Figure 10–3.

Figure 10–3

Scatter Plot for Example 10–3



After the plot is drawn, it should be analyzed to determine which type of relationship, if any, exists. For example, the plot shown in Figure 10–1 suggests a positive relationship, since as the number of cars rented increases, revenue tends to increase also. The plot of the data shown in Figure 10–2 suggests a negative relationship, since as the number of absences increases, the final grade decreases. Finally, the plot of the data shown in Figure 10–3 shows no specific type of relationship, since no pattern is discernible.

Note that the data shown in Figures 10–1 and 10–2 also suggest a linear relationship, since the points seem to fit a straight line, although not perfectly. Sometimes a scatter plot, such as the one in Figure 10–4, shows a curvilinear relationship between the data. In this situation, the methods shown in this section and in Section 10–2 cannot be used. Methods for curvilinear relationships are beyond the scope of this book.

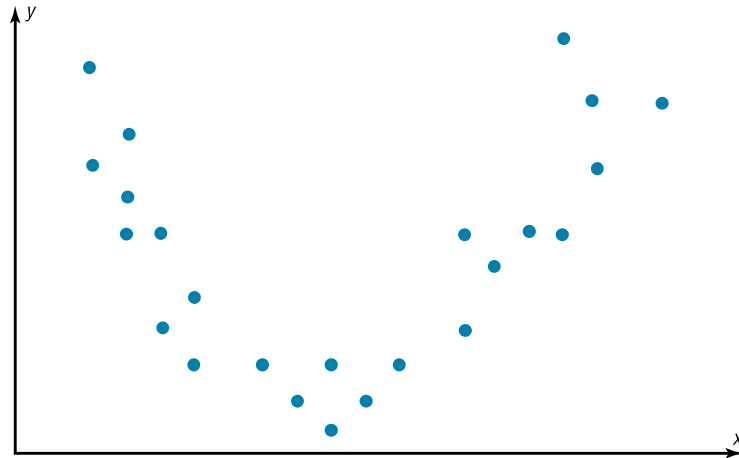
Correlation**Objective 2**

Compute the correlation coefficient.

Correlation Coefficient As stated in the Introduction, statisticians use a measure called the *correlation coefficient* to determine the strength of the linear relationship between two variables. There are several types of correlation coefficients. The one

Figure 10–4

Scatter Plot
Suggesting a
Curvilinear
Relationship



explained in this section is called the **Pearson product moment correlation coefficient (PPMC)**, named after statistician Karl Pearson, who pioneered the research in this area.

The **correlation coefficient** computed from the sample data measures the strength and direction of a linear relationship between two quantitative variables. The symbol for the sample correlation coefficient is r . The symbol for the population correlation coefficient is ρ (Greek letter rho).

The *range of the correlation coefficient* is from -1 to $+1$. If there is a *strong positive linear relationship* between the variables, the value of r will be close to $+1$. If there is a *strong negative linear relationship* between the variables, the value of r will be close to -1 . When there is no linear relationship between the variables or only a weak relationship, the value of r will be close to 0 . See Figure 10–5.

The graphs in Figure 10–6 show the relationship between the correlation coefficients and their corresponding scatter plots. Notice that as the value of the correlation coefficient increases from 0 to $+1$ (parts *a*, *b*, and *c*), data values become closer to an increasingly strong relationship. As the value of the correlation coefficient decreases from 0 to -1 (parts *d*, *e*, and *f*), the data values also become closer to a straight line. Again this suggests a stronger relationship.

There are several ways to compute the value of the correlation coefficient. One method is to use the formula shown here.

Formula for the Correlation Coefficient r

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

where n is the number of data pairs.

Figure 10–5

Range of Values for the
Correlation Coefficient

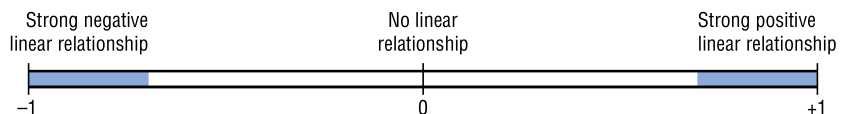
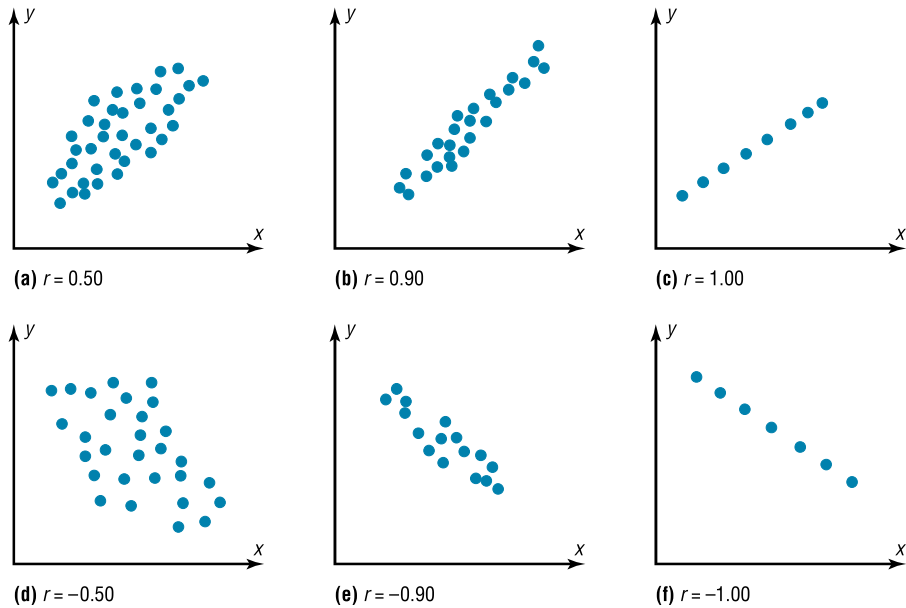


Figure 10–6

Relationship Between the Correlation Coefficient and the Scatter Plot



Assumptions for the Correlation Coefficient

1. The sample is a random sample.
2. The data pairs fall approximately on a straight line and are measured at the interval or ratio level.
3. The variables have a joint normal distribution. (This means that given any specific value of x , the y values are normally distributed; and given any specific value of y , the x values are normally distributed.)

Rounding Rule for the Correlation Coefficient Round the value of r to three decimal places.

The formula looks somewhat complicated, but using a table to compute the values, as shown in Example 10–4, makes it somewhat easier to determine the value of r .

There are no units associated with r , and the value of r will remain unchanged if the x and y values are switched.

Example 10–4

Car Rental Companies

Compute the correlation coefficient for the data in Example 10–1.

Solution

Step 1 Make a table as shown here.

Company	Cars x (in ten thousands)	Revenue y (in billions)	xy	x^2	y^2
A	63.0	7.0			
B	29.0	3.9			
C	20.8	2.1			
D	19.1	2.8			
E	13.4	1.4			
F	8.5	1.5			

Step 2 Find the values of xy , x^2 , and y^2 and place these values in the corresponding columns of the table.

The completed table is shown.

Company	Cars x (in 10,000s)	Revenue y (in billions)	xy	x^2	y^2
A	63.0	7.0	441.00	3969.00	49.00
B	29.0	3.9	113.10	841.00	15.21
C	20.8	2.1	43.68	432.64	4.41
D	19.1	2.8	53.48	364.81	7.84
E	13.4	1.4	18.76	179.56	1.96
F	8.5	1.5	12.75	72.25	2.25
	$\Sigma x = 153.8$	$\Sigma y = 18.7$	$\Sigma xy = 682.77$	$\Sigma x^2 = 5859.26$	$\Sigma y^2 = 80.67$

Step 3 Substitute in the formula and solve for r :

$$\begin{aligned}
 r &= \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n(\Sigma x^2) - (\Sigma x)^2][n(\Sigma y^2) - (\Sigma y)^2]}} \\
 &= \frac{(6)(682.77) - (153.8)(18.7)}{\sqrt{[(6)(5859.26) - (153.8)^2][(6)(80.67) - (18.7)^2]}} = 0.982
 \end{aligned}$$

The correlation coefficient suggests a strong relationship between the number of cars a rental agency has and its annual revenue.

Example 10–5

Absences and Final Grades

Compute the value of the correlation coefficient for the data obtained in the study of the number of absences and the final grade of the seven students in the statistics class given in Example 10–2.

Solution

Step 1 Make a table.

Step 2 Find the values of xy , x^2 , and y^2 ; place these values in the corresponding columns of the table.

Student	Number of absences x	Final grade y (%)	xy	x^2	y^2
A	6	82	492	36	6,724
B	2	86	172	4	7,396
C	15	43	645	225	1,849
D	9	74	666	81	5,476
E	12	58	696	144	3,364
F	5	90	450	25	8,100
G	8	78	624	64	6,084
	$\Sigma x = 57$	$\Sigma y = 511$	$\Sigma xy = 3745$	$\Sigma x^2 = 579$	$\Sigma y^2 = 38,993$

Step 3 Substitute in the formula and solve for r :

$$\begin{aligned}
 r &= \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n(\Sigma x^2) - (\Sigma x)^2][n(\Sigma y^2) - (\Sigma y)^2]}} \\
 &= \frac{(7)(3745) - (57)(511)}{\sqrt{[(7)(579) - (57)^2][(7)(38,993) - (511)^2]}} = -0.944
 \end{aligned}$$

The value of r suggests a strong negative relationship between a student's final grade and the number of absences a student has. That is, the more absences a student has, the lower is his or her grade.

Example 10–6

Age and Wealth

Compute the value of the correlation coefficient for the data given in Example 10–3 for the age and wealth of the richest persons in the United States.

Solution

Step 1 Make a table.

Step 2 Find the values of xy , x^2 , and y^2 , and place these values in the corresponding columns of the table.

Person	Age x	Net wealth y	xy	x^2	y^2
A	73	16	1,168	5,329	256
B	65	26	1,690	4,225	676
C	53	50	2,650	2,809	2,500
D	54	21.5	1,161	2,916	462.25
E	79	40	3,160	6,241	1,600
F	69	16	1,104	4,761	256
G	61	19.6	1,195.6	3,721	384.16
H	65	19	1,235	4,225	361
	$\Sigma x = 519$	$\Sigma y = 208.1$	$\Sigma xy = 13,363.6$	$\Sigma x^2 = 34,227$	$\Sigma y^2 = 6,495.41$

Step 3 Substitute in the formula and solve for r :

$$\begin{aligned}
 r &= \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n(\Sigma x^2) - (\Sigma x)^2][n(\Sigma y^2) - (\Sigma y)^2]}} \\
 &= \frac{8(13,363.6) - (519)(208.1)}{\sqrt{[8(34,227) - (519)^2][8(6,495.41) - (208.1)^2]}} \\
 &= \frac{-1095.1}{\sqrt{(4455)(8657.67)}} \\
 &= \frac{-1095.1}{6210.469} \\
 &= -0.176
 \end{aligned}$$

The value of r indicates a very weak negative relationship between the variables.

In Example 10–4, the value of r was high (close to 1.00); in Example 10–6, the value of r was much lower (close to 0). This question then arises, When is the value of r due to chance, and when does it suggest a significant linear relationship between the variables? This question will be answered next.

Objective 3

Test the hypothesis $H_0: \rho = 0$.

The Significance of the Correlation Coefficient As stated before, the range of the correlation coefficient is between -1 and $+1$. When the value of r is near $+1$ or -1 , there is a strong linear relationship. When the value of r is near 0 , the linear relationship is weak or nonexistent. Since the value of r is computed from data obtained from samples, there are two possibilities when r is not equal to zero: either the value of r is high enough to conclude that there is a significant linear relationship between the variables, or the value of r is due to chance.

To make this decision, you use a hypothesis-testing procedure. The traditional method is similar to the one used in previous chapters.

- Step 1** State the hypotheses.
- Step 2** Find the critical values.
- Step 3** Compute the test value.
- Step 4** Make the decision.
- Step 5** Summarize the results.

The population correlation coefficient is computed from taking all possible (x, y) pairs; it is designated by the Greek letter ρ (rho). The sample correlation coefficient can then be used as an estimator of ρ if the following assumptions are valid.

1. The variables x and y are *linearly* related.
2. The variables are *random* variables.
3. The two variables have a *bivariate normal distribution*.

A bivariate normal distribution means that for the pairs of (x, y) data values, the corresponding y values have a bell-shaped distribution for any given x value, and the x values for any given y value have a bell-shaped distribution.

Formally defined, the **population correlation coefficient** ρ is the correlation computed by using all possible pairs of data values (x, y) taken from a population.

Interesting Fact

Scientists think that a person is never more than 3 feet away from a spider at any given time!

In hypothesis testing, one of these is true:

- $H_0: \rho = 0$ This null hypothesis means that there is no correlation between the x and y variables in the population.
- $H_1: \rho \neq 0$ This alternative hypothesis means that there is a significant correlation between the variables in the population.

When the null hypothesis is rejected at a specific level, it means that there is a significant difference between the value of r and 0. When the null hypothesis is not rejected, it means that the value of r is not significantly different from 0 (zero) and is probably due to chance.

Several methods can be used to test the significance of the correlation coefficient. Three methods will be shown in this section. The first uses the t test.

Historical Notes

A mathematician named Karl Pearson (1857–1936) became interested in Francis Galton's work and saw that the correlation and regression theory could be applied to other areas besides heredity. Pearson developed the correlation coefficient that bears his name.

Formula for the t Test for the Correlation Coefficient

$$t = r\sqrt{\frac{n-2}{1-r^2}}$$

with degrees of freedom equal to $n - 2$.

Although hypothesis tests can be one-tailed, most hypotheses involving the correlation coefficient are two-tailed. Recall that ρ represents the population correlation coefficient. Also, if there is no linear relationship, the value of the correlation coefficient will be 0. Hence, the hypotheses will be

$$H_0: \rho = 0 \quad \text{and} \quad H_1: \rho \neq 0$$

You do not have to identify the claim here, since the question will always be whether there is a significant linear relationship between the variables.

The two-tailed critical values are used. These values are found in Table F in Appendix C. Also, when you are testing the significance of a correlation coefficient, both variables x and y must come from normally distributed populations.

Example 10-7

Test the significance of the correlation coefficient found in Example 10-4. Use $\alpha = 0.05$ and $r = 0.982$.

Solution

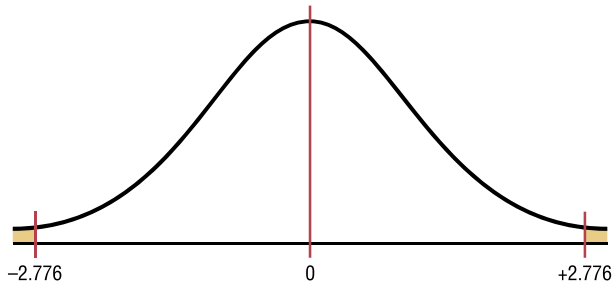
Step 1 State the hypotheses.

$$H_0: \rho = 0 \quad \text{and} \quad H_1: \rho \neq 0$$

Step 2 Find the critical values. Since $\alpha = 0.05$ and there are $6 - 2 = 4$ degrees of freedom, the critical values obtained from Table F are ± 2.776 , as shown in Figure 10-7.

Figure 10-7

Critical Values for Example 10-7



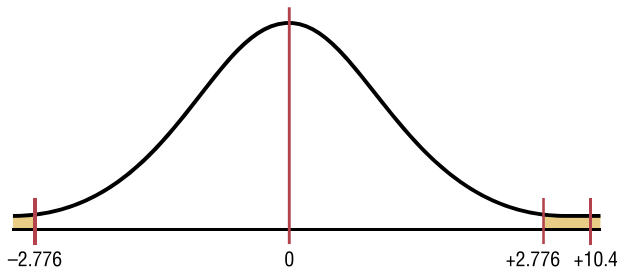
Step 3 Compute the test value.

$$t = r \sqrt{\frac{n - 2}{1 - r^2}} = 0.982 \sqrt{\frac{6 - 2}{1 - (0.982)^2}} = 10.4$$

Step 4 Make the decision. Reject the null hypothesis, since the test value falls in the critical region, as shown in Figure 10-8.

Figure 10-8

Test Value for Example 10-7



Step 5 Summarize the results. There is a significant relationship between the number of cars a rental agency owns and its annual income.

The second method that can be used to test the significance of r is the P -value method. The method is the same as that shown in Chapters 8 and 9. It uses the following steps.

Step 1 State the hypotheses.

Step 2 Find the test value. (In this case, use the t test.)

- Step 3** Find the P -value. (In this case, use Table F.)
- Step 4** Make the decision.
- Step 5** Summarize the results.

Consider an example where $t = 4.059$ and $d.f. = 4$. Using Table F with $d.f. = 4$ and the row Two tails, the value 4.059 falls between 3.747 and 4.604; hence, $0.01 < P\text{-value} < 0.02$. (The P -value obtained from a calculator is 0.015.) That is, the P -value falls between 0.01 and 0.02. The decision, then, is to reject the null hypothesis since $P\text{-value} < 0.05$.

The third method of testing the significance of r is to use Table I in Appendix C. This table shows the values of the correlation coefficient that are significant for a specific α level and a specific number of degrees of freedom. For example, for 7 degrees of freedom and $\alpha = 0.05$, the table gives a critical value of 0.666. Any value of r greater than +0.666 or less than -0.666 will be significant, and the null hypothesis will be rejected. See Figure 10–9. When Table I is used, you need not compute the t test value. Table I is for two-tailed tests only.

Figure 10–9
Finding the Critical Value from Table I

d.f.	$\alpha = 0.05$	$\alpha = 0.01$
1		
2		
3		
4		
5		
6		
7	0.666	

Example 10–8

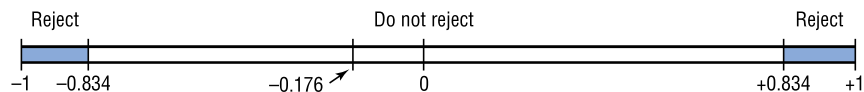
Using Table I, test the significance at $\alpha = 0.01$ of the correlation coefficient $r = -0.176$, obtained in Example 10–6.

Solution

$$H_0: \rho = 0 \quad \text{and} \quad H_1: \rho \neq 0$$

Since the sample size is 8, there are $n - 2$, or $8 - 2 = 6$, degrees of freedom. When $\alpha = 0.01$ and $d.f. = 6$, the value obtained from Table I is 0.834. For a significant relationship, a value of r greater than +0.834 or less than -0.834 is needed. Since the value of $r = -0.176$ is greater than -0.834 , the null hypothesis is not rejected. Hence there is not enough evidence to say that there is a significant linear relationship between age and wealth. See Figure 10–10.

Figure 10–10
Rejection and Nonrejection Regions for Example 10–8





Correlation and Causation Researchers must understand the nature of the linear relationship between the independent variable x and the dependent variable y . When a hypothesis test indicates that a significant linear relationship exists between the variables, researchers must consider the possibilities outlined next.

Possible Relationships Between Variables

When the null hypothesis has been rejected for a specific α value, any of the following five possibilities can exist.

1. *There is a direct cause-and-effect relationship between the variables.* That is, x causes y . For example, water causes plants to grow, poison causes death, and heat causes ice to melt.
2. *There is a reverse cause-and-effect relationship between the variables.* That is, y causes x . For example, suppose a researcher believes excessive coffee consumption causes nervousness, but the researcher fails to consider that the reverse situation may occur. That is, it may be that an extremely nervous person craves coffee to calm his or her nerves.
3. *The relationship between the variables may be caused by a third variable.* For example, if a statistician correlated the number of deaths due to drowning and the number of cans of soft drink consumed daily during the summer, he or she would probably find a significant relationship. However, the soft drink is not necessarily responsible for the deaths, since both variables may be related to heat and humidity.
4. *There may be a complexity of interrelationships among many variables.* For example, a researcher may find a significant relationship between students' high school grades and college grades. But there probably are many other variables involved, such as IQ, hours of study, influence of parents, motivation, age, and instructors.
5. *The relationship may be coincidental.* For example, a researcher may be able to find a significant relationship between the increase in the number of people who are exercising and the increase in the number of people who are committing crimes. But common sense dictates that any relationship between these two values must be due to coincidence.

When two variables are highly correlated, item 3 in the box states that there exists a possibility that the correlation is due to a third variable. If this is the case and the third variable is unknown to the researcher or not accounted for in the study, it is called a

lurking variable. An attempt should be made by the researcher to identify such variables and to use methods to control their influence.

It is important to restate the fact that even if the correlation between two variables is high, it does not necessarily mean causation. There are other possibilities, such as lurking variables or just a coincidental relationship. See the Speaking of Statistics article on page 548.

Also, you should be cautious when the data for one or both of the variables involve averages rather than individual data. It is not wrong to use averages, but the results cannot be generalized to individuals since averaging tends to smooth out the variability among individual data values. The result could be a higher correlation than actually exists.

Thus, when the null hypothesis is rejected, the researcher must consider all possibilities and select the appropriate one as determined by the study. Remember, correlation does not necessarily imply causation.

Applying the Concepts 10–1

Stopping Distances

In a study on speed control, it was found that the main reasons for regulations were to make traffic flow more efficient and to minimize the risk of danger. An area that was focused on in the study was the distance required to completely stop a vehicle at various speeds. Use the following table to answer the questions.

MPH	Braking distance (feet)
20	20
30	45
40	81
50	133
60	205
80	411

Assume MPH is going to be used to predict stopping distance.

- Which of the two variables is the independent variable?
- Which is the dependent variable?
- What type of variable is the independent variable?
- What type of variable is the dependent variable?
- Construct a scatter plot for the data.
- Is there a linear relationship between the two variables?
- Redraw the scatter plot, and change the distances between the independent-variable numbers. Does the relationship look different?
- Is the relationship positive or negative?
- Can braking distance be accurately predicted from MPH?
- List some other variables that affect braking distance.
- Compute the value of r .
- Is r significant at $\alpha = 0.05$?

See page 589 for the answers.

Extending the Concepts

28. One of the formulas for computing r is

$$r = \frac{\sum(x - \bar{x})(y - \bar{y})}{(n - 1)(s_x)(s_y)}$$

Using the data in Exercise 27, compute r with this formula. Compare the results.

29. Compute r for the data set shown. Explain the reason for this value of r . Now, interchange the values of x and y and compute r again. Compare this value with the previous one. Explain the results of the comparison.

x	1	2	3	4	5
y	3	5	7	9	11

30. Compute r for the following data and test the hypothesis $H_0: \rho = 0$. Draw the scatter plot; then explain the results.

x	-3	-2	-1	0	1	2	3
y	9	4	1	0	1	4	9

10-2

Regression

Objective 4

Compute the equation of the regression line.

In studying relationships between two variables, collect the data and then construct a scatter plot. The purpose of the scatter plot, as indicated previously, is to determine the nature of the relationship. The possibilities include a positive linear relationship, a negative linear relationship, a curvilinear relationship, or no discernible relationship. After the scatter plot is drawn, the next steps are to compute the value of the correlation coefficient and to test the significance of the relationship. If the value of the correlation coefficient is significant, the next step is to determine the equation of the **regression line**, which is the data's line of best fit. (*Note:* Determining the regression line when r is not significant and then making predictions using the regression line are meaningless.) The purpose of the regression line is to enable the researcher to see the trend and make predictions on the basis of the data.

Line of Best Fit

Figure 10-11 shows a scatter plot for the data of two variables. It shows that several lines can be drawn on the graph near the points. Given a scatter plot, you must be able to draw the *line of best fit*. *Best fit* means that the sum of the squares of the vertical distances from

Figure 10-11

Scatter Plot with Three Lines Fit to the Data

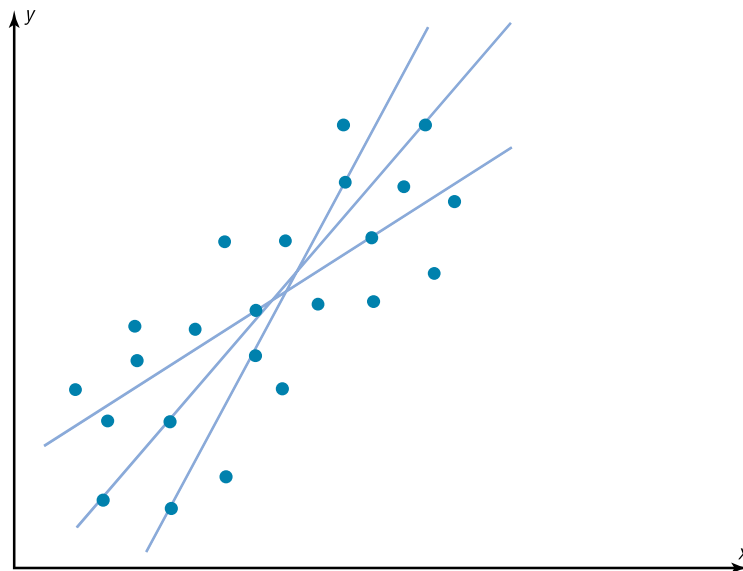
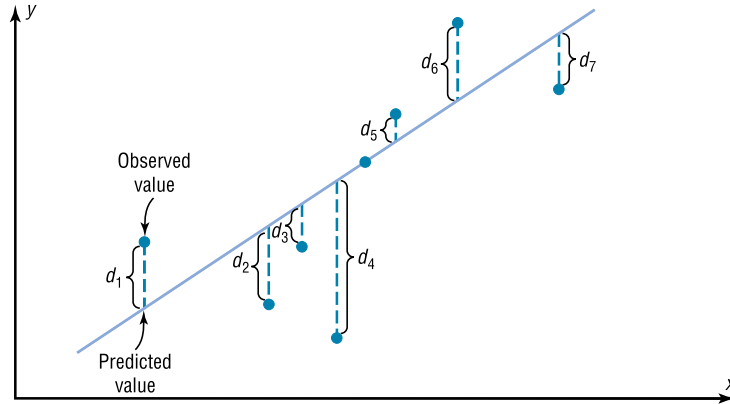


Figure 10–12
Line of Best Fit for a Set of Data Points



Historical Notes
Francis Galton drew the line of best fit visually. An assistant of Karl Pearson's named G. Yule devised the mathematical solution using the least-squares method, employing a mathematical technique developed by Adrien-Marie Legendre about 100 years earlier.

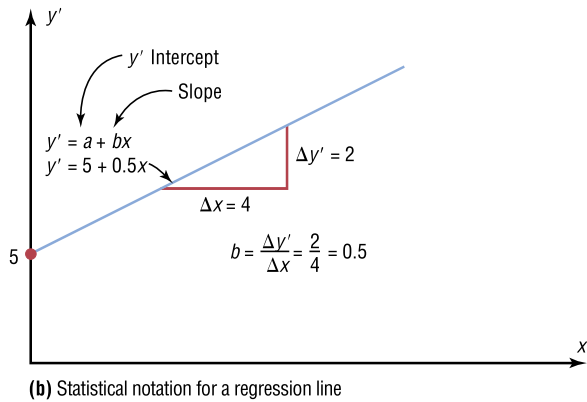
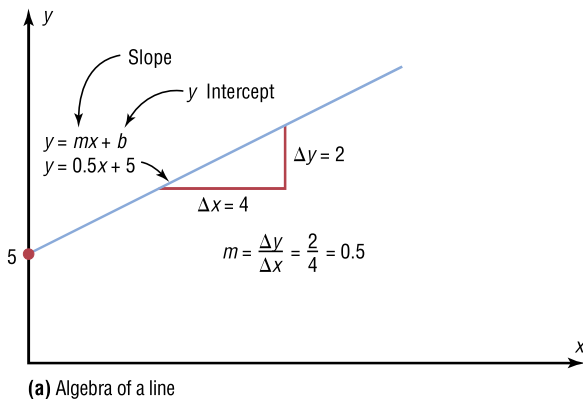
each point to the line is at a minimum. The reason you need a line of best fit is that the values of y will be predicted from the values of x ; hence, the closer the points are to the line, the better the fit and the prediction will be. See Figure 10–12. When r is positive, the line slopes upward and to the right. When r is negative, the line slopes downward from left to right.

Determination of the Regression Line Equation

In algebra, the equation of a line is usually given as $y = mx + b$, where m is the slope of the line and b is the y intercept. (Students who need an algebraic review of the properties of a line should refer to Appendix A, Section A–3, before studying this section.) In statistics, the equation of the regression line is written as $y' = a + bx$, where a is the y' intercept and b is the slope of the line. See Figure 10–13.

There are several methods for finding the equation of the regression line. Two formulas are given here. *These formulas use the same values that are used in computing the value of the correlation coefficient.* The mathematical development of these formulas is beyond the scope of this book.

Figure 10–13
A Line as Represented in Algebra and in Statistics



Formulas for the Regression Line $y' = a + bx$

$$a = \frac{(\Sigma y)(\Sigma x^2) - (\Sigma x)(\Sigma xy)}{n(\Sigma x^2) - (\Sigma x)^2}$$

$$b = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2}$$

where a is the y' intercept and b is the slope of the line.

Rounding Rule for the Intercept and Slope Round the values of a and b to three decimal places.

Example 10–9**Car Rental Companies**

Find the equation of the regression line for the data in Example 10–4, and graph the line on the scatter plot of the data.

Solution

The values needed for the equation are $n = 6$, $\Sigma x = 153.8$, $\Sigma y = 18.7$, $\Sigma xy = 682.77$, and $\Sigma x^2 = 5859.26$. Substituting in the formulas, you get

$$a = \frac{(\Sigma y)(\Sigma x^2) - (\Sigma x)(\Sigma xy)}{n(\Sigma x^2) - (\Sigma x)^2} = \frac{(18.7)(5859.26) - (153.8)(682.77)}{(6)(5859.26) - (153.8)^2} = 0.396$$

$$b = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2} = \frac{6(682.77) - (153.8)(18.7)}{(6)(5859.26) - (153.8)^2} = 0.106$$

Hence, the equation of the regression line $y' = a + bx$ is

$$y' = 0.396 + 0.106x$$

To graph the line, select any two points for x and find the corresponding values for y . Use any x values between 10 and 60. For example, let $x = 15$. Substitute in the equation and find the corresponding y' value.

$$\begin{aligned} y' &= 0.396 \\ &= 0.396 + 0.106(15) \\ &= 1.986 \end{aligned}$$

Let $x = 40$; then

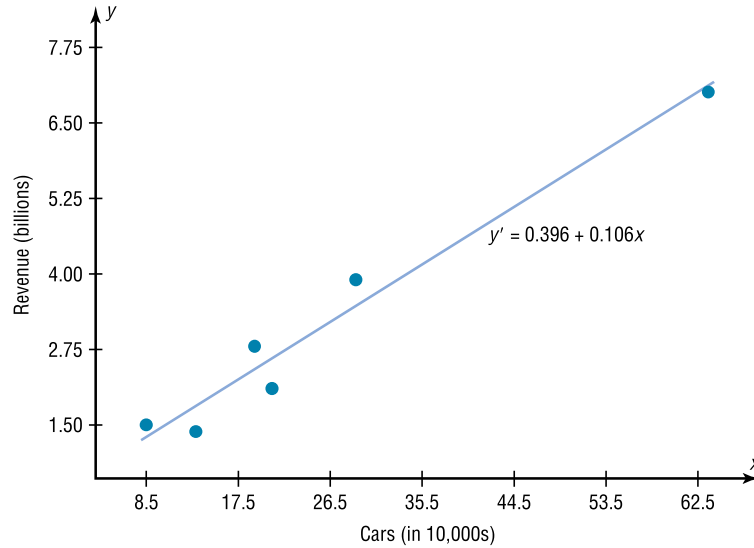
$$\begin{aligned} y' &= 0.396 + 0.106x \\ &= 0.396 + 0.106(40) \\ &= 4.636 \end{aligned}$$

Then plot the two points (15, 1.986) and (40, 4.636) and draw a line connecting the two points. See Figure 10–14.

Note: When you draw the regression line, it is sometimes necessary to *truncate* the graph (see Chapter 2). This is done when the distance between the origin and the first labeled coordinate on the x axis is not the same as the distance between the rest of the

Figure 10–14

Regression Line for Example 10–9



labeled x coordinates or the distance between the origin and the first labeled y' coordinate is not the same as the distance between the other labeled y' coordinates. When the x axis or the y axis has been truncated, do not use the y' intercept value to graph the line. When you graph the regression line, always select x values between the smallest x data value and the largest x data value.

Example 10–10

Absences and Final Grades

Find the equation of the regression line for the data in Example 10–5, and graph the line on the scatter plot.

Solution

The values needed for the equation are $n = 7$, $\Sigma x = 57$, $\Sigma y = 511$, $\Sigma xy = 3745$, and $\Sigma x^2 = 579$. Substituting in the formulas, you get

$$a = \frac{(\Sigma y)(\Sigma x^2) - (\Sigma x)(\Sigma xy)}{n(\Sigma x^2) - (\Sigma x)^2} = \frac{(511)(579) - (57)(3745)}{(7)(579) - (57)^2} = 102.493$$

$$b = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2} = \frac{(7)(3745) - (57)(511)}{(7)(579) - (57)^2} = -3.622$$

Hence, the equation of the regression line $y' = a + bx$ is

$$y' = 102.493 - 3.622x$$

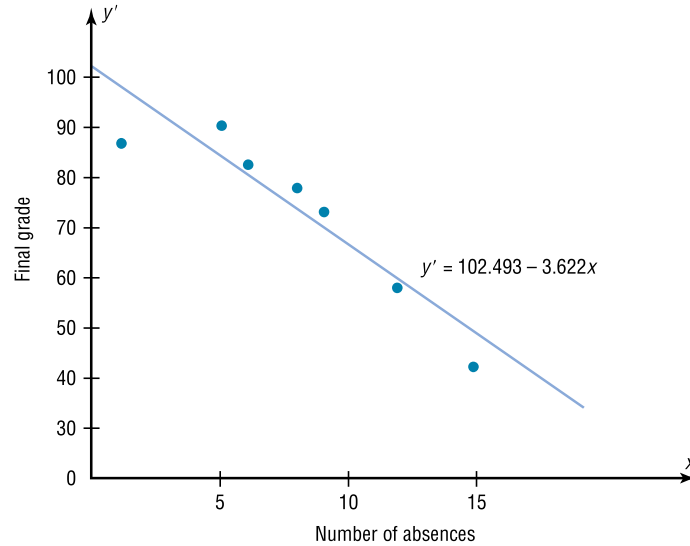
The graph of the line is shown in Figure 10–15.

Historical Note

In 1795, Adrien-Marie Legendre (1752–1833) measured the meridian arc on the earth's surface from Barcelona, Spain, to Dunkirk, England. This measure was used as the basis for the measure of the meter. Legendre developed the least-squares method around the year 1805.

The sign of the correlation coefficient and the sign of the slope of the regression line will always be the same. That is, if r is positive, then b will be positive; if r is negative, then b will be negative. The reason is that the numerators of the formulas are the same and determine the signs of r and b , and the denominators are always positive. The regression line will always pass through the point whose x coordinate is the mean of the x values and whose y coordinate is the mean of the y values, that is, (\bar{x}, \bar{y}) .

Figure 10–15
Regression Line for
Example 10–10



The regression line can be used to make predictions for the dependent variable. The method for making predictions is shown in Example 10–11.

Example 10–11

Car Rental Companies

Use the equation of the regression line to predict the income of a car rental agency that has 200,000 automobiles.

Solution

Since the x values are in 10,000s, divide 200,000 by 10,000 to get 20, and then substitute 20 for x in the equation.

$$\begin{aligned} y' &= 0.396 + 0.106x \\ &= 0.396 + 0.106(20) \\ &= 2.516 \end{aligned}$$

Hence, when a rental agency has 200,000 automobiles, its revenue will be approximately \$2.516 billion.

The value obtained in Example 10–11 is a point prediction, and with point predictions, no degree of accuracy or confidence can be determined. More information on prediction is given in Section 10–3.

The magnitude of the change in one variable when the other variable changes exactly 1 unit is called a **marginal change**. The value of slope b of the regression line equation represents the marginal change. For example, in Example 10–9 the slope of the regression line is 0.106, which means for each increase of 10,000 cars, the value of y changes 0.106 unit (\$106 million) on average.

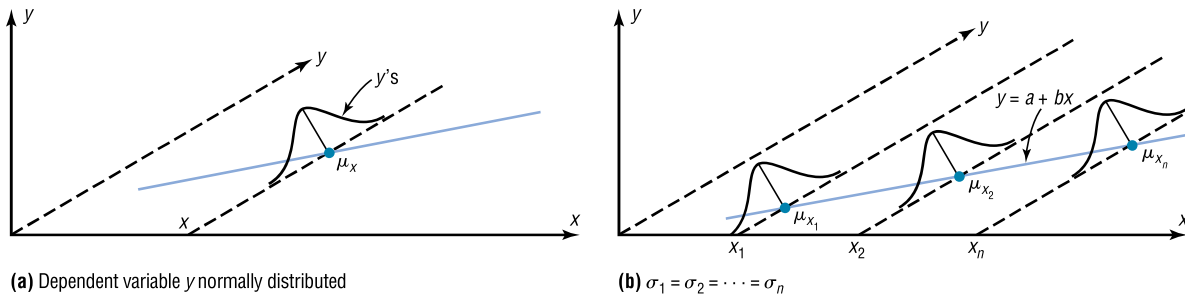
When r is not significantly different from 0, the best predictor of y is the mean of the data values of y . For valid predictions, the value of the correlation coefficient must be significant. Also, two other assumptions must be met.

Assumptions for Valid Predictions in Regression

1. The sample is a random sample.
2. For any specific value of the independent variable x , the value of the dependent variable y must be normally distributed about the regression line. See Figure 10–16(a).
3. The standard deviation of each of the dependent variables must be the same for each value of the independent variable. See Figure 10–16(b).

Figure 10–16

Assumptions for Predictions



Extrapolation, or making predictions beyond the bounds of the data, must be interpreted cautiously. For example, in 1979, some experts predicted that the United States would run out of oil by the year 2003. This prediction was based on the current consumption and on known oil reserves at that time. However, since then, the automobile industry has produced many new fuel-efficient vehicles. Also, there are many as yet undiscovered oil fields. Finally, science may someday discover a way to run a car on something as unlikely but as common as peanut oil. In addition, the price of a gallon of gasoline was predicted to reach \$10 a few years later. Fortunately this has not come to pass. *Remember that when predictions are made, they are based on present conditions or on the premise that present trends will continue.* This assumption may or may not prove true in the future.

The steps for finding the value of the correlation coefficient and the regression line equation are summarized in this Procedure Table:

Interesting Fact

It is estimated that wearing a motorcycle helmet reduces the risk of a fatal accident by 30%.

Procedure Table				
Finding the Correlation Coefficient and the Regression Line Equation				
Step 1	Make a table, as shown in step 2.			
Step 2	Find the values of xy , x^2 , and y^2 . Place them in the appropriate columns and sum each column.			
x	y	xy	x^2	y^2
·	·	·	·	·
·	·	·	·	·
·	·	·	·	·
$\Sigma x =$ 	$\Sigma y =$ 	$\Sigma xy =$ 	$\Sigma x^2 =$ 	$\Sigma y^2 =$

Procedure Table (Continued)**Step 3** Substitute in the formula to find the value of r .

$$r = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n(\Sigma x^2) - (\Sigma x)^2][n(\Sigma y^2) - (\Sigma y)^2]}}$$

Step 4 When r is significant, substitute in the formulas to find the values of a and b for the regression line equation $y' = a + bx$.

$$a = \frac{(\Sigma y)(\Sigma x^2) - (\Sigma x)(\Sigma xy)}{n(\Sigma x^2) - (\Sigma x)^2} \quad b = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2}$$

A scatter plot should be checked for outliers. An outlier is a point that seems out of place when compared with the other points (see Chapter 3). Some of these points can affect the equation of the regression line. When this happens, the points are called **influential points** or **influential observations**.

When a point on the scatter plot appears to be an outlier, it should be checked to see if it is an influential point. An influential point tends to “pull” the regression line toward the point itself. To check for an influential point, the regression line should be graphed with the point included in the data set. Then a second regression line should be graphed that excludes the point from the data set. If the position of the second line is changed considerably, the point is said to be an influential point. Points that are outliers in the x direction tend to be influential points.

Researchers should use their judgment as to whether to include influential observations in the final analysis of the data. If the researcher feels that the observation is not necessary, then it should be excluded so that it does not influence the results of the study. However, if the researcher feels that it is necessary, then he or she may want to obtain additional data values whose x values are near the x value of the influential point and then include them in the study.



© Dave Carpenter. King Features Syndicate.