

## Lesson

## 2-2

## Linear Models

## Vocabulary

linear function  
 linear model  
 interpolation  
 extrapolation  
 observed values  
 predicted values  
 residual  
 sum of squared residuals

► **BIG IDEA** The sum of squared deviations is a statistic for determining which of two lines fits the data better.

A **linear function** is a set of ordered pairs  $(x, y)$  satisfying an equation of the form  $y = mx + b$ , where  $m$  and  $b$  are constants. Recall that the graph of every such function is a line with slope  $m$  and  $y$ -intercept  $b$ .

## Fitting a Line to Data

When data in a scatterplot lie near a line, we can create a **linear model** for the data, that is, a model of one variable as a linear function of the other. Even if the linear function does not contain all of the data points, it may still be useful in describing the overall trend of the data or in predicting values of either the dependent or independent variable.

In this lesson and the next, you will learn several techniques for constructing linear models. One technique is to fit a line to data “by eye,” that is, to draw a line that seems close to the data.

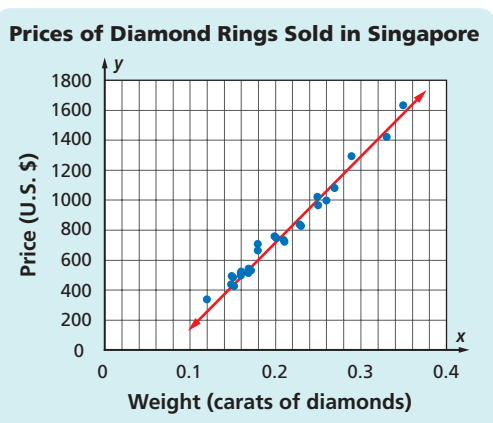
## Mental Math

What is the sum of a set of 25 elements whose mean is 300?

## GUIDED

## Example 1

Jewelers emphasize that the price of a diamond is determined by cut, carat weight, color and clarity. The table at the right gives carat weights and approximate prices in U.S. dollars for twenty diamond rings sold at a recent auction in Singapore. All rings are of the same quality gold and contain a single diamond. The data are graphed below.



Weight $x$	Price $y$ (U.S. \$)
0.18	702.00
0.17	517.50
0.25	963.00
0.29	1290.00
0.27	1080.00
0.15	484.50
0.20	747.00
0.25	1017.00
0.21	724.50
0.17	529.50
0.35	1629.00
0.33	1417.50
0.26	994.50
0.16	513.00
0.12	334.50
0.18	664.50
0.15	430.50
0.16	507.00
0.16	498.00
0.23	829.50

Source: Journal of Statistics Education

The linear model is based on the weight of the diamond used. Although the data are not collinear, the line through (0.18, 600) and (0.32, 1400) seems close to the points. It has been added to the graph. An equation of this graphed line is one linear model for these data. Is the size alone a good predictor of price?

- Find an equation of the graphed line which relates weight and price.
- Interpret the slope of the line in the context of the problem.
- Use the model to estimate the price of a 0.3-carat diamond ring.
- Why is the model not good for predicting the cost of a 0.05-carat diamond ring?
- Why is the set of data not a function?

### Solution

- We use the points (0.18, 600) and (0.32, 1400) to find an equation of the line:

$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{?}{?} \approx ?$$

Substitute  $m = 5714.29$  and (0.18, 600) into the point-slope form of the equation for a line.

$$\begin{aligned} y - y_1 &= m(x - x_1) \\ y - \frac{?}{?} &= 5714.29(x - \frac{?}{?}) \\ y &= \frac{?}{?} \end{aligned}$$

- The slope is the rate of change. The cost increases by about  $\frac{?}{?}$  dollars for every 1-carat increase in weight.
- To predict the cost from the weight, substitute  $x = 0.3$  and solve for  $y$ .  
When  $x = 0.3$ ,  $y = \frac{?}{?}$ . According to the model, the cost of a 0.3-carat diamond ring will be about  $\frac{?}{?}$ .
- Substitute 0.05 for  $x$  and solve for  $y$ . The model predicts that a 0.05-carat diamond ring will cost  $\frac{?}{?}$ . This means that the seller would pay you to take the ring! This is not plausible, so the model is not a good predictor.
- There is at least one  $x$ -value that does not have a unique  $y$ -value. There are two diamonds at 0.17 carats that have different prices, so the data set is not a function.



### Hard Rock

The largest rough diamond ever found weighed 3106 carats.

In this example, known diamond weights range from 0.12 to 0.35 carats. If you use the model to predict the price of a 0.3-carat diamond, you are making a prediction *between* known values. Prediction between known values is called **interpolation**. If you calculate the price of a diamond weighing 0.05 carats, you are making a prediction *beyond* known values. Prediction beyond known values is called **extrapolation**. Extrapolation is usually more hazardous than interpolation, because it depends on an assumption that a relationship will continue past the known data. In this case, a diamond much smaller than those in the sample might be inexpensive, but it will not be free.

## Measuring How Well a Line Models Data

In Chapter 1, you studied the sample variance, which was computed as the sum of the squared deviations from the mean divided by  $n - 1$ . A similar statistic, the *residual*, tells how far away data are from your chosen model. At the right is a table showing a smaller sample of the diamond ring data on page 87. These data, collected from sources such as experiments or surveys, are called **observed values**. Below, the scatterplot of these data is shown together with the graph of the line with equation  $y = 2400x + 400$ . This equation seems to be a pretty good model for the data. However, is it the best model? To compare models, we calculate *residuals*. The values predicted by a model are called **predicted values**. The observed value minus the predicted value is the **residual**. A residual is positive when the observed value is higher than what is predicted by the model. A residual is negative when the observed value is lower than what the model predicts. For instance, a 0.16-carat diamond ring sold for \$507.00. This is lower than the price predicted by the model.

$$\begin{aligned}\text{predicted price} &= 2400 \cdot (0.16) + 400 \\ &= \$784\end{aligned}$$

The residual, or error, is

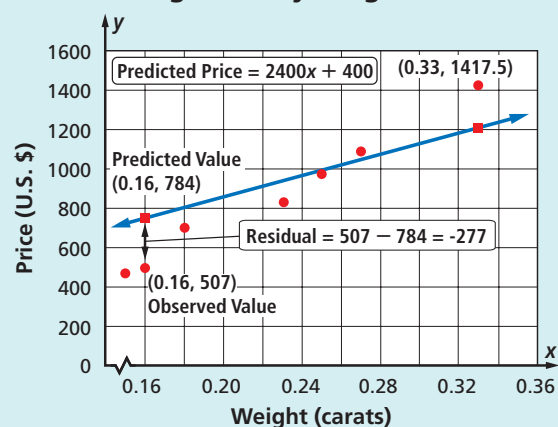
$$\begin{aligned}\text{residual} &= \text{observed value} - \text{predicted value} \\ &= 507 - 784 \\ &= -277.\end{aligned}$$

The actual cost was \$277 less than predicted. On the graph, we see that the observed value is 277 units below the linear model  $y = 2400x + 400$ . That is why the residual is negative. Every data point has a residual, so  $n$  data points provide  $n$  residuals. The absolute value of a residual is the length of the vertical segment from the data point to the corresponding point on the linear model.

Diamond Ring Prices by Weight of Diamond

Weight	Price (U.S. dollars)
0.15	484.50
0.16	507.00
0.18	702.00
0.25	963.00
0.27	1080.00
0.33	1417.50
0.23	829.50

Diamond Ring Prices by Weight of Diamond



### ► QY

What is the residual for the 0.33-carat diamond?

On the graphs below, this line and another linear model are drawn to fit the seven points. Recall that when you compute variance, you add the squared deviations from the mean. Similarly, when modeling data with a line, we can measure the variation from the line by adding the squares of the residuals. You can use spreadsheets to calculate residuals and the sum of squared residuals, as shown below.

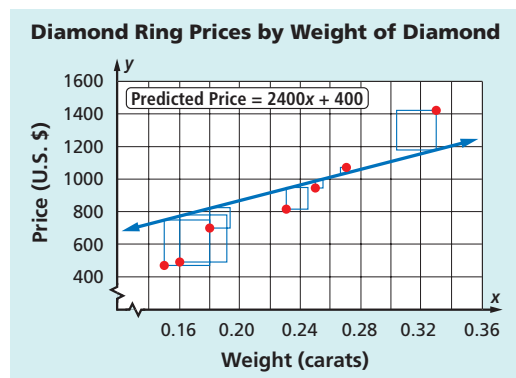
	A weight	B price	C predicted	D residuals
			=2400*a[]+400	=b[]-c[]
1	0.15	484.5	760.	-275.5
2	0.16	507	784.	-277.
3	0.18	702	832.	-130.
4	0.25	963	1000.	-37.
5	0.27	1080	1048.	32.
D	residuals:=b[]-c[]			

	E predicted	D residuals	E sqresid	F sumsq...
	=2400*a[]+400	=b[]-c[]	=d[]^2	
1	760.	-275.5	75900.3	237779.
2	784.	-277.	76729.	
3	832.	-130.	16900.	
4	1000.	-37.	1369.	
5	1048.	32.	1024.	
F1	=sum(e[])			

You can also think of  $(\text{residual})^2$  as the area of a geometric square whose side has length equal to the absolute value of the residual. These geometric squares are drawn on the graphs below.

### Linear Model 1

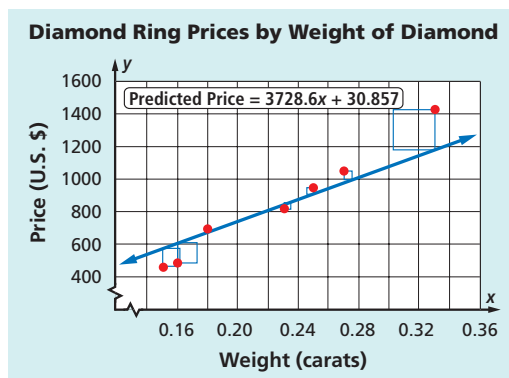
Squares are shown for a line that does not go through any data points.



Total area of the squares  $\approx 237,800$

### Linear Model 2

Squares are shown for a line that goes through two of the data points.



Total area of the squares  $\approx 59,870$

The second line is a better model of the data because it has a smaller total area of the squares. The total area is the **sum of squared residuals**.

### Definition of Sum of Squared Residuals

$$\text{Sum of squared residuals} = \sum_{i=1}^n (\text{observed } y_i - \text{predicted } y_i)^2$$

The sum of squared residuals is a statistic that measures lack of fit. If you compare two lines, the one with the larger sum of squared residuals is not as good a model as the one with the smaller sum of squared residuals. If you have many possible lines you could use to model data, compute the sum of squared residuals for each model. The line that gives the smallest value provides the best fit to the data.

## Activity

This table shows the number of televisions per 100 people in 1997 and the number of unemployed per 100 people in 2008 for nine countries from around the world. Some people suggest that increased TV viewing leads to a less productive workforce. Are these statistics related?

**Step 1** Use a statistics utility to create a scatterplot with number of TVs on the x-axis and number of unemployed people on the y-axis.

**Step 2** Add a movable line to the scatterplot.

**Step 3** Move the line as close to  $y = -0.3x + 17$  as you can. Your screen should look similar to the one at the right. Record the sum of squared residuals.

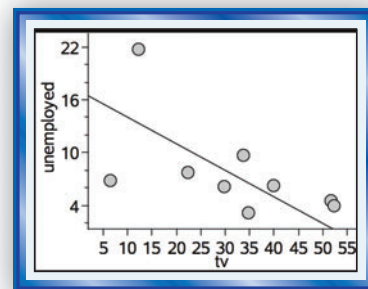
**Step 4** Move the line so that it goes through the points for Bulgaria and the Netherlands. Record the equation. Record the sum of squared residuals.

**Step 5** Tell which line is the better model of the data. Explain your choice.

**Step 6** Move the line until you get the smallest sum of squares that you can. Record the equation and the sum of squared residuals.

Country	TVs per 100	Unemployed per 100
Argentina	22.3	7.8
Bulgaria	40.0	6.3
India	6.5	6.8
Israel	29.9	6.1
Netherlands	51.8	4.5
New Zealand	52.3	4.0
Poland	33.7	9.7
South Africa	12.3	21.7
South Korea	34.7	3.2

Source: UNESCO Institute for Statistics, CIA World Fact Book



## Questions

### COVERING THE IDEAS

- Define *linear function*.
- Refer to the diamond price data in Example 1.
  - Find an equation for the line that passes through the data points (0.35, 1629) and (0.12, 334.50).
  - Write a sentence that states how cost increases for every carat increase in weight according to your equation.
  - Use your line to predict the price of a 0.17-carat diamond.
  - Is this interpolation or extrapolation?
- Fill in the Blanks** For a data set, residuals are the differences between \_\_\_?\_\_\_ values and \_\_\_?\_\_\_ values.

4. When a residual is positive, is the observed value higher or lower than the predicted value?

In 5 and 6, a diamond speculator used the line with equation  $y = 5000x - 250$  to estimate the price of diamond rings.

5. a. What would the speculator predict for the price of the 0.29-carat diamond ring?  
b. According to the data in Example 1, what is the residual for the 0.29-carat diamond ring?
6. a. What would the speculator predict as a price for the 0.21-carat diamond ring?  
b. What is the residual for this diamond ring?
7. Does the phrase “sum of squared residuals” mean “first sum the residuals and then square the sum” or “first square each residual and then sum the squares?”
8. Consider the linear model  $f(x) = 5.2x - 3$ .  
a. Given a residual of  $-0.2$  at  $x = 1$ , find the observed value.  
b. Given a residual of  $2.1$  at  $x = 4$ , find the observed value.  
c. Given a residual of  $0$  at  $x = 3$ , find the observed value.

### APPLYING THE MATHEMATICS

In 9 and 10, suppose Jane asked members of her family how many states they had visited. Her data are in the table at the right.

9. Jane plotted the data and drew a line through the points for Ed and Grandma. She found its equation to be  $y = 0.55x + 3.5$ .  
a. What is the slope of this line and what does it represent?  
b. Complete the spreadsheet below to calculate the residuals for this model.  
c. What is the sum of squared residuals for Jane’s model?  
d. Jane used her equation to predict how many states she will have visited at age 30. Is this interpolation or extrapolation?  
e. What is Jane’s prediction from Part d?

Family Member	Age	States Visited
Jane	16	15
Cousin Xia	16	18
Dad	40	27
Grandma	70	42
Brother Ed	10	9
Uncle Ralph	45	19

◇	A	B	C	D
1	Age	States Visited	Predicted	Residual <sup>2</sup>
2	16	15	12.3	7.29
3	16	18		
4	40	27	25.5	2.25
5	70	42	42.0	0
6	10	9		
7	45	19	28.25	85.56

10. Xia thought that drawing the line through the points for Ed and Dad would look better because then there would be two points above the line and two below the line. Her line has equation  $y = 0.6x + 3$ .



