



Practice Exercises: Lesson 6.2

Diez, D. M., Çetinkaya-Rundel, M., Barr, C. D. (2019). OpenIntro Statistics (4th ed.). OpenIntro.
<https://www.openintro.org/book/os/> CC BY-SA 3.0

STAT 1201
Introduction to Probability and Statistics

ONLINE AND DISTANCE EDUCATION

Exercises

8.17 Units of regression. Consider a regression predicting weight (kg) from height (cm) for a sample of adult males. What are the units of the correlation coefficient, the intercept, and the slope?

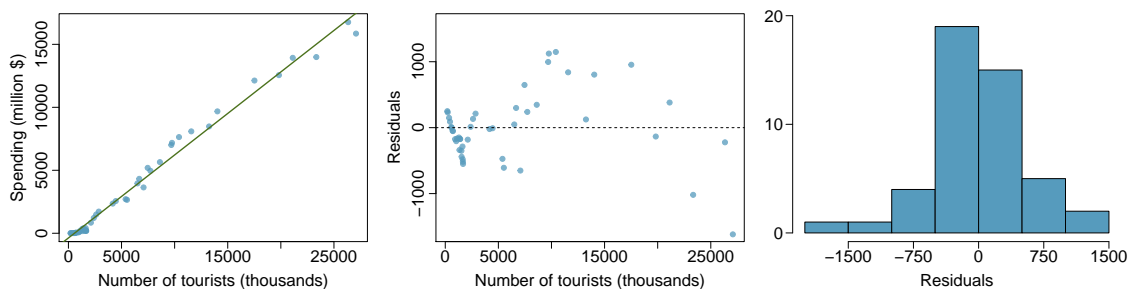
8.18 Which is higher? Determine if I or II is higher or if they are equal. Explain your reasoning. For a regression line, the uncertainty associated with the slope estimate, b_1 , is higher when

- I. there is a lot of scatter around the regression line or
- II. there is very little scatter around the regression line

8.19 Over-under, Part I. Suppose we fit a regression line to predict the shelf life of an apple based on its weight. For a particular apple, we predict the shelf life to be 4.6 days. The apple's residual is -0.6 days. Did we over or under estimate the shelf-life of the apple? Explain your reasoning.

8.20 Over-under, Part II. Suppose we fit a regression line to predict the number of incidents of skin cancer per 1,000 people from the number of sunny days in a year. For a particular year, we predict the incidence of skin cancer to be 1.5 per 1,000 people, and the residual for this year is 0.5. Did we over or under estimate the incidence of skin cancer? Explain your reasoning.

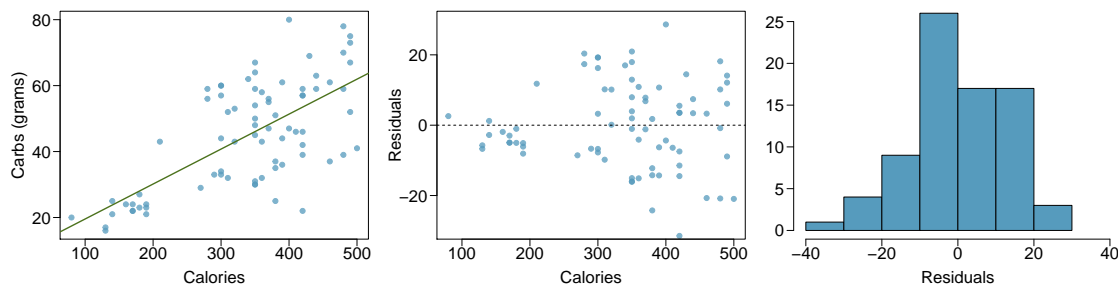
8.21 Tourism spending. The Association of Turkish Travel Agencies reports the number of foreign tourists visiting Turkey and tourist spending by year.¹⁴ Three plots are provided: scatterplot showing the relationship between these two variables along with the least squares fit, residuals plot, and histogram of residuals.



- (a) Describe the relationship between number of tourists and spending.
- (b) What are the explanatory and response variables?
- (c) Why might we want to fit a regression line to these data?
- (d) Do the data meet the conditions required for fitting a least squares line? In addition to the scatterplot, use the residual plot and histogram to answer this question.

¹⁴Association of Turkish Travel Agencies, Foreign Visitors Figure & Tourist Spendings By Years.

8.22 Nutrition at Starbucks, Part I. The scatterplot below shows the relationship between the number of calories and amount of carbohydrates (in grams) Starbucks food menu items contain.¹⁵ Since Starbucks only lists the number of calories on the display items, we are interested in predicting the amount of carbs a menu item has based on its calorie content.



- Describe the relationship between number of calories and amount of carbohydrates (in grams) that Starbucks food menu items contain.
- In this scenario, what are the explanatory and response variables?
- Why might we want to fit a regression line to these data?
- Do these data meet the conditions required for fitting a least squares line?

8.23 The Coast Starlight, Part II. Exercise 8.11 introduces data on the Coast Starlight Amtrak train that runs from Seattle to Los Angeles. The mean travel time from one stop to the next on the Coast Starlight is 129 mins, with a standard deviation of 113 minutes. The mean distance traveled from one stop to the next is 108 miles with a standard deviation of 99 miles. The correlation between travel time and distance is 0.636.

- Write the equation of the regression line for predicting travel time.
- Interpret the slope and the intercept in this context.
- Calculate R^2 of the regression line for predicting travel time from distance traveled for the Coast Starlight, and interpret R^2 in the context of the application.
- The distance between Santa Barbara and Los Angeles is 103 miles. Use the model to estimate the time it takes for the Starlight to travel between these two cities.
- It actually takes the Coast Starlight about 168 mins to travel from Santa Barbara to Los Angeles. Calculate the residual and explain the meaning of this residual value.
- Suppose Amtrak is considering adding a stop to the Coast Starlight 500 miles away from Los Angeles. Would it be appropriate to use this linear model to predict the travel time from Los Angeles to this point?

8.24 Body measurements, Part III. Exercise 8.13 introduces data on shoulder girth and height of a group of individuals. The mean shoulder girth is 107.20 cm with a standard deviation of 10.37 cm. The mean height is 171.14 cm with a standard deviation of 9.41 cm. The correlation between height and shoulder girth is 0.67.

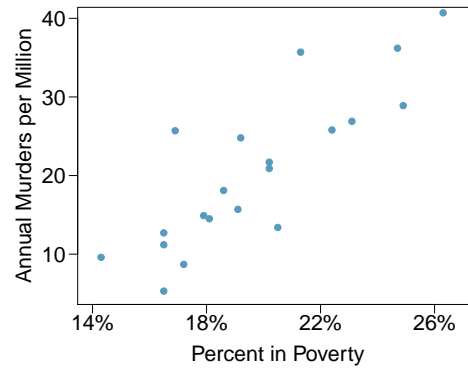
- Write the equation of the regression line for predicting height.
- Interpret the slope and the intercept in this context.
- Calculate R^2 of the regression line for predicting height from shoulder girth, and interpret it in the context of the application.
- A randomly selected student from your class has a shoulder girth of 100 cm. Predict the height of this student using the model.
- The student from part (d) is 160 cm tall. Calculate the residual, and explain what this residual means.
- A one year old has a shoulder girth of 56 cm. Would it be appropriate to use this linear model to predict the height of this child?

¹⁵Source: Starbucks.com, collected on March 10, 2011, www.starbucks.com/menu/nutrition.

8.25 Murders and poverty, Part I. The following regression output is for predicting annual murders per million from percentage living in poverty in a random sample of 20 metropolitan areas.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-29.901	7.789	-3.839	0.001
poverty%	2.559	0.390	6.562	0.000
$s = 5.512$	$R^2 = 70.52\%$		$R^2_{adj} = 68.89\%$	

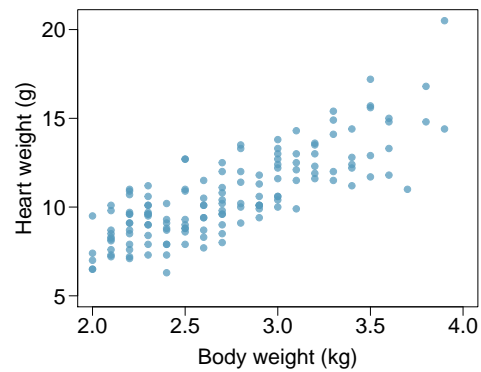
- Write out the linear model.
- Interpret the intercept.
- Interpret the slope.
- Interpret R^2 .
- Calculate the correlation coefficient.



8.26 Cats, Part I. The following regression output is for predicting the heart weight (in g) of cats from their body weight (in kg). The coefficients are estimated using a dataset of 144 domestic cats.

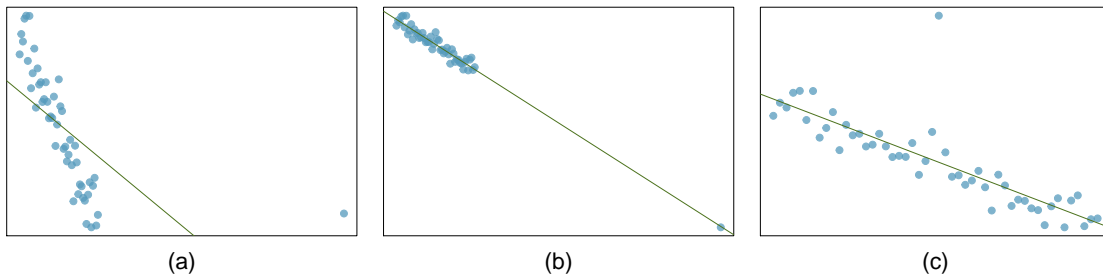
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.357	0.692	-0.515	0.607
body wt	4.034	0.250	16.119	0.000
$s = 1.452$	$R^2 = 64.66\%$		$R^2_{adj} = 64.41\%$	

- Write out the linear model.
- Interpret the intercept.
- Interpret the slope.
- Interpret R^2 .
- Calculate the correlation coefficient.

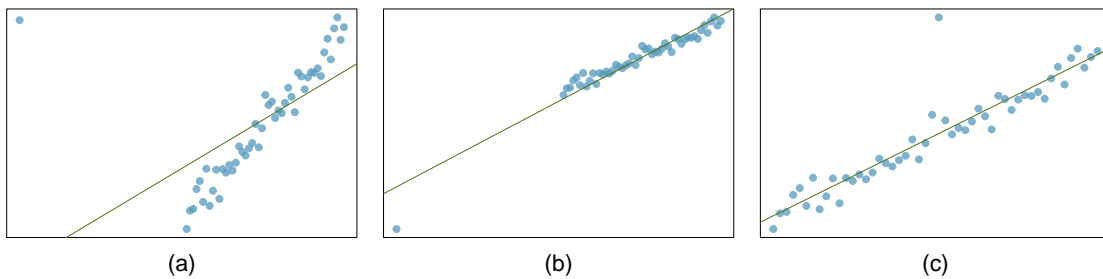


Exercises

8.27 Outliers, Part I. Identify the outliers in the scatterplots shown below, and determine what type of outliers they are. Explain your reasoning.

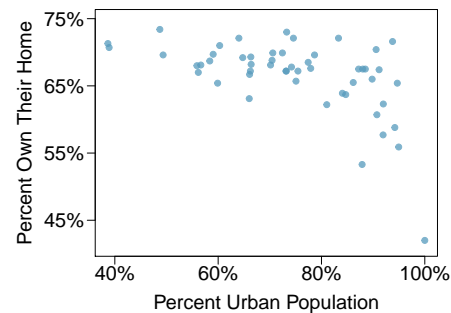


8.28 Outliers, Part II. Identify the outliers in the scatterplots shown below and determine what type of outliers they are. Explain your reasoning.



8.29 Urban homeowners, Part I. The scatterplot below shows the percent of families who own their home vs. the percent of the population living in urban areas.¹⁶ There are 52 observations, each corresponding to a state in the US. Puerto Rico and District of Columbia are also included.

- Describe the relationship between the percent of families who own their home and the percent of the population living in urban areas.
- The outlier at the bottom right corner is District of Columbia, where 100% of the population is considered urban. What type of an outlier is this observation?



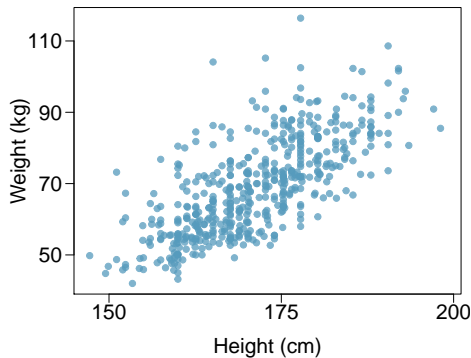
8.30 Crawling babies, Part II. Exercise 8.12 introduces data on the average monthly temperature during the month babies first try to crawl (about 6 months after birth) and the average first crawling age for babies born in a given month. A scatterplot of these two variables reveals a potential outlying month when the average temperature is about 53°F and average crawling age is about 28.5 weeks. Does this point have high leverage? Is it an influential point?

¹⁶United States Census Bureau, 2010 Census Urban and Rural Classification and Urban Area Criteria and Housing Characteristics: 2010.

Exercises

In the following exercises, visually check the conditions for fitting a least squares regression line. However, you do not need to report these conditions in your solutions.

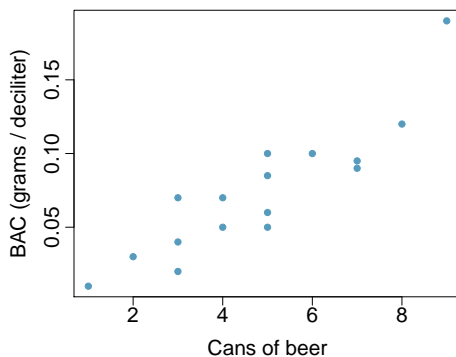
8.31 Body measurements, Part IV. The scatterplot and least squares summary below show the relationship between weight measured in kilograms and height measured in centimeters of 507 physically active individuals.



	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-105.0113	7.5394	-13.93	0.0000
height	1.0176	0.0440	23.13	0.0000

- Describe the relationship between height and weight.
- Write the equation of the regression line. Interpret the slope and intercept in context.
- Do the data provide strong evidence that an increase in height is associated with an increase in weight? State the null and alternative hypotheses, report the p-value, and state your conclusion.
- The correlation coefficient for height and weight is 0.72. Calculate R^2 and interpret it in context.

8.32 Beer and blood alcohol content. Many people believe that gender, weight, drinking habits, and many other factors are much more important in predicting blood alcohol content (BAC) than simply considering the number of drinks a person consumed. Here we examine data from sixteen student volunteers at Ohio State University who each drank a randomly assigned number of cans of beer. These students were evenly divided between men and women, and they differed in weight and drinking habits. Thirty minutes later, a police officer measured their blood alcohol content (BAC) in grams of alcohol per deciliter of blood.¹⁹ The scatterplot and regression table summarize the findings.

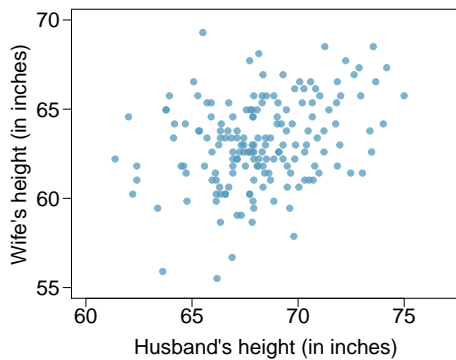


	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.0127	0.0126	-1.00	0.3320
beers	0.0180	0.0024	7.48	0.0000

- Describe the relationship between the number of cans of beer and BAC.
- Write the equation of the regression line. Interpret the slope and intercept in context.
- Do the data provide strong evidence that drinking more cans of beer is associated with an increase in blood alcohol? State the null and alternative hypotheses, report the p-value, and state your conclusion.
- The correlation coefficient for number of cans of beer and BAC is 0.89. Calculate R^2 and interpret it in context.
- Suppose we visit a bar, ask people how many drinks they have had, and also take their BAC. Do you think the relationship between number of drinks and BAC would be as strong as the relationship found in the Ohio State study?

¹⁹J. Malkevitch and L.M. Lesser. *For All Practical Purposes: Mathematical Literacy in Today's World*. WH Freeman & Co, 2008.

8.33 Husbands and wives, Part II. The scatterplot below summarizes husbands' and wives' heights in a random sample of 170 married couples in Britain, where both partners' ages are below 65 years. Summary output of the least squares fit for predicting wife's height from husband's height is also provided in the table.

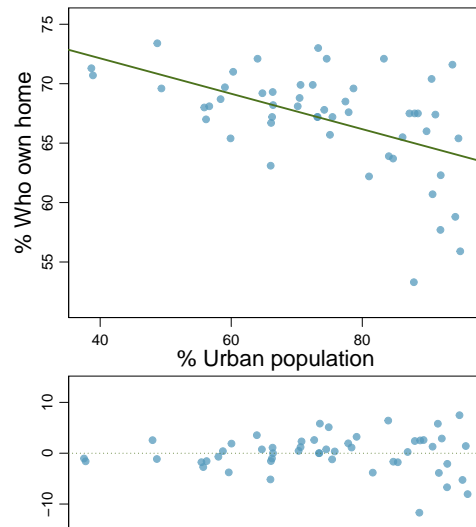


	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	43.5755	4.6842	9.30	0.0000
height_husband	0.2863	0.0686	4.17	0.0000

- Is there strong evidence that taller men marry taller women? State the hypotheses and include any information used to conduct the test.
- Write the equation of the regression line for predicting wife's height from husband's height.
- Interpret the slope and intercept in the context of the application.
- Given that $R^2 = 0.09$, what is the correlation of heights in this data set?
- You meet a married man from Britain who is 5'9" (69 inches). What would you predict his wife's height to be? How reliable is this prediction?
- You meet another married man from Britain who is 6'7" (79 inches). Would it be wise to use the same linear model to predict his wife's height? Why or why not?

8.34 Urban homeowners, Part II. Exercise 8.29 gives a scatterplot displaying the relationship between the percent of families that own their home and the percent of the population living in urban areas. Below is a similar scatterplot, excluding District of Columbia, as well as the residuals plot. There were 51 cases.

- For these data, $R^2 = 0.28$. What is the correlation? How can you tell if it is positive or negative?
- Examine the residual plot. What do you observe? Is a simple least squares fit appropriate for these data?

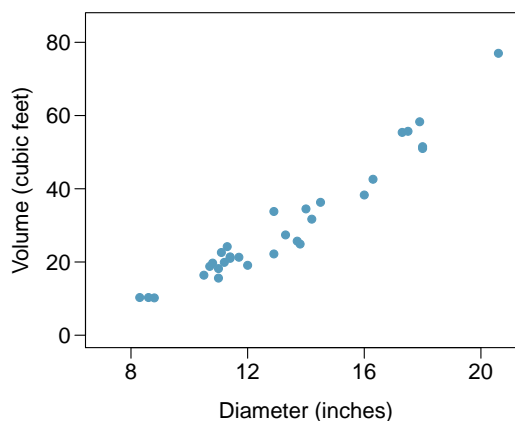
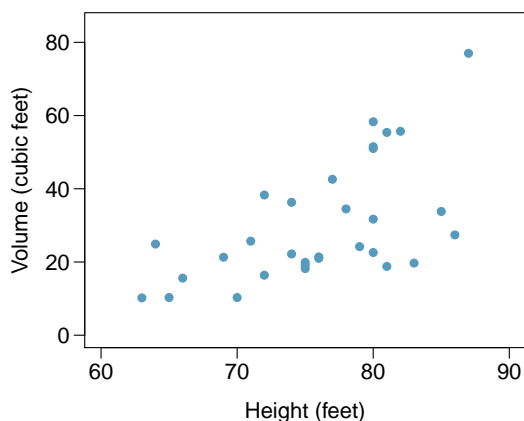


Chapter exercises

8.37 True / False. Determine if the following statements are true or false. If false, explain why.

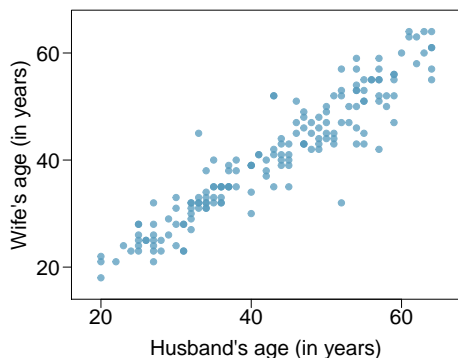
- A correlation coefficient of -0.90 indicates a stronger linear relationship than a correlation of 0.5 .
- Correlation is a measure of the association between any two variables.

8.38 Trees. The scatterplots below show the relationship between height, diameter, and volume of timber in 31 felled black cherry trees. The diameter of the tree is measured 4.5 feet above the ground.²⁰



- Describe the relationship between volume and height of these trees.
- Describe the relationship between volume and diameter of these trees.
- Suppose you have height and diameter measurements for another black cherry tree. Which of these variables would be preferable to use to predict the volume of timber in this tree using a simple linear regression model? Explain your reasoning.

8.39 Husbands and wives, Part III. Exercise 8.33 presents a scatterplot displaying the relationship between husbands' and wives' ages in a random sample of 170 married couples in Britain, where both partners' ages are below 65 years. Given below is summary output of the least squares fit for predicting wife's age from husband's age.



	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.5740	1.1501	1.37	0.1730
age_husband	0.9112	0.0259	35.25	0.0000

$df = 168$

- We might wonder, is the age difference between husbands and wives consistent across ages? If this were the case, then the slope parameter would be $\beta_1 = 1$. Use the information above to evaluate if there is strong evidence that the difference in husband and wife ages differs for different ages.
- Write the equation of the regression line for predicting wife's age from husband's age.
- Interpret the slope and intercept in context.
- Given that $R^2 = 0.88$, what is the correlation of ages in this data set?
- You meet a married man from Britain who is 55 years old. What would you predict his wife's age to be? How reliable is this prediction?
- You meet another married man from Britain who is 85 years old. Would it be wise to use the same linear model to predict his wife's age? Explain.

²⁰Source: R Dataset, stat.ethz.ch/R-manual/R-patched/library/datasets/html/trees.html.

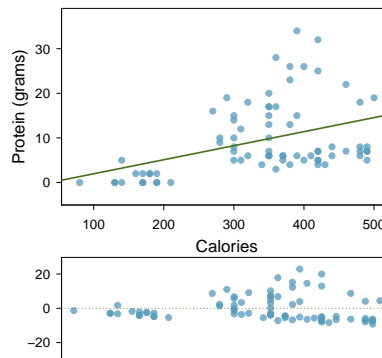
8.40 Cats, Part II. Exercise 8.26 presents regression output from a model for predicting the heart weight (in g) of cats from their body weight (in kg). The coefficients are estimated using a dataset of 144 domestic cat. The model output is also provided below.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.357	0.692	-0.515	0.607
body wt	4.034	0.250	16.119	0.000

$$s = 1.452 \quad R^2 = 64.66\% \quad R_{adj}^2 = 64.41\%$$

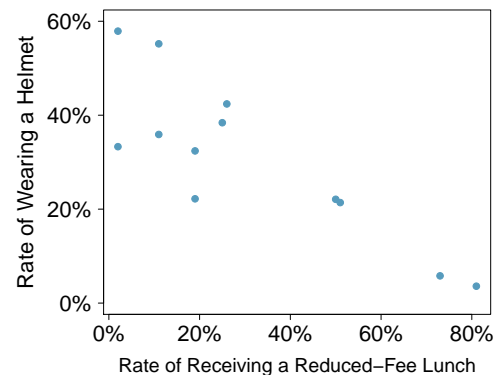
- We see that the point estimate for the slope is positive. What are the hypotheses for evaluating whether body weight is positively associated with heart weight in cats?
- State the conclusion of the hypothesis test from part (a) in context of the data.
- Calculate a 95% confidence interval for the slope of body weight, and interpret it in context of the data.
- Do your results from the hypothesis test and the confidence interval agree? Explain.

8.41 Nutrition at Starbucks, Part II. Exercise 8.22 introduced a data set on nutrition information on Starbucks food menu items. Based on the scatterplot and the residual plot provided, describe the relationship between the protein content and calories of these menu items, and determine if a simple linear model is appropriate to predict amount of protein from the number of calories.



8.42 Helmets and lunches. The scatterplot shows the relationship between socioeconomic status measured as the percentage of children in a neighborhood receiving reduced-fee lunches at school (**lunch**) and the percentage of bike riders in the neighborhood wearing helmets (**helmet**). The average percentage of children receiving reduced-fee lunches is 30.8% with a standard deviation of 26.7% and the average percentage of bike riders wearing helmets is 38.8% with a standard deviation of 16.9%.

- If the R^2 for the least-squares regression line for these data is 72%, what is the correlation between **lunch** and **helmet**?
- Calculate the slope and intercept for the least-squares regression line for these data.
- Interpret the intercept of the least-squares regression line in the context of the application.
- Interpret the slope of the least-squares regression line in the context of the application.
- What would the value of the residual be for a neighborhood where 40% of the children receive reduced-fee lunches and 40% of the bike riders wear helmets? Interpret the meaning of this residual in the context of the application.



8.43 Match the correlation, Part III. Match each correlation to the corresponding scatterplot.

- $r = -0.72$
- $r = 0.07$
- $r = 0.86$
- $r = 0.99$

