

Chapter 6: Exploring Data: Relationships

Lesson Plan

- Displaying Relationships: Scatterplots
- Making Predictions: Regression Line
- Correlation
- Least-Squares Regression
- Interpreting Correlation and Regression

For All Practical
Purposes



Mathematical Literacy in
Today's World, 9th ed.

Chapter 6: Exploring Data: Distributions

Displaying Relationships

■ Relationship Between Two Variables

- Examine data for two variables to see if there is a relationship between the variables. Does one influence the other?
- Study both variables on the same individual.
- If a relationship exists between variables, typically one variable influences or causes a change in another variable.
 - Explanatory variable explains, or causes, the change in another variable.
 - Response variable measures the outcome, or response to the change.

Response variable –
A variable that measures an outcome or result of a study (observed outcome).

Explanatory variable –
A variable that explains or causes change in the response variable.

Chapter 6: Exploring Data: Distributions

Displaying Relationships: Scatterplots

■ Data to Be Used for a Scatterplot

- A scatterplot is a graph that shows the relationship between two numerical variables, measured on the same individual.
 - Explanatory variable, x , is plotted on the horizontal axis (x).
 - Response variable, y , is plotted on the vertical axis (y).
 - Each pair of related variables (x, y) is plotted on the graph.

Example: A study is done to see how the number of beers that a student drinks predicts his/her blood alcohol content (BAC). *Results of 16 students:*

Student	1	2	3	4	5	6	7	8
Beers	5	2	9	8	3	7	3	5
BAC	0.10	0.03	0.19	0.12	0.04	0.095	0.07	0.06
Student	9	10	11	12	13	14	15	16
Beers	3	5	4	6	5	7	1	4
BAC	0.02	0.05	0.07	0.10	0.85	0.09	0.01	0.05

Explanatory variable,
 $x =$ beers drunk

Response variable,
 $y =$ BAC level

Chapter 6: Exploring Data: Distributions

Displaying Relationships: Scatterplots

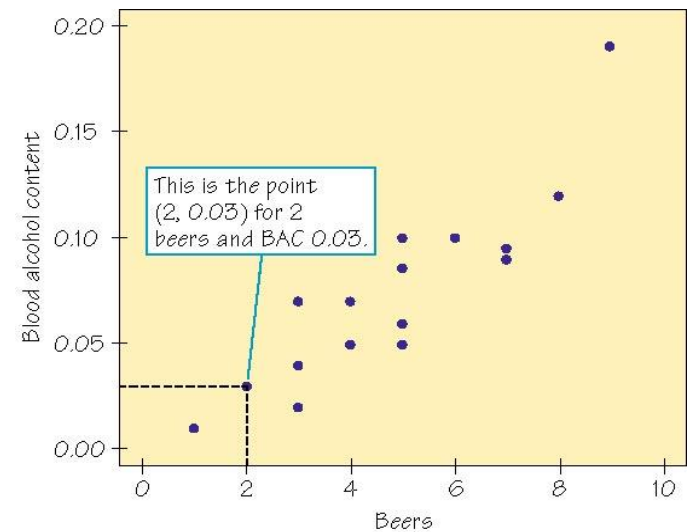
■ Scatterplot

- Example continued: The scatterplot of the blood alcohol content, **BAC**, (**y**, response variable) against the number of beers a young adult drinks (**x**, explanatory variable).
- The data from the previous table are plotted as points on the graph (**x**, **y**).

Examining This Scatterplot...

1. What is the overall pattern (form, direction, and strength)?
 - Form – Roughly a straight-line pattern.
 - Direction – Positive association (both increase).
 - Strength – Moderately strong (mostly on line).
2. Any striking deviations (outliers)? Not here.

BAC vs. number of beers consumed



Outliers – A deviation in a distribution of a data point falling outside the overall pattern.

Chapter 6: Exploring Data: Distributions

Regression Lines

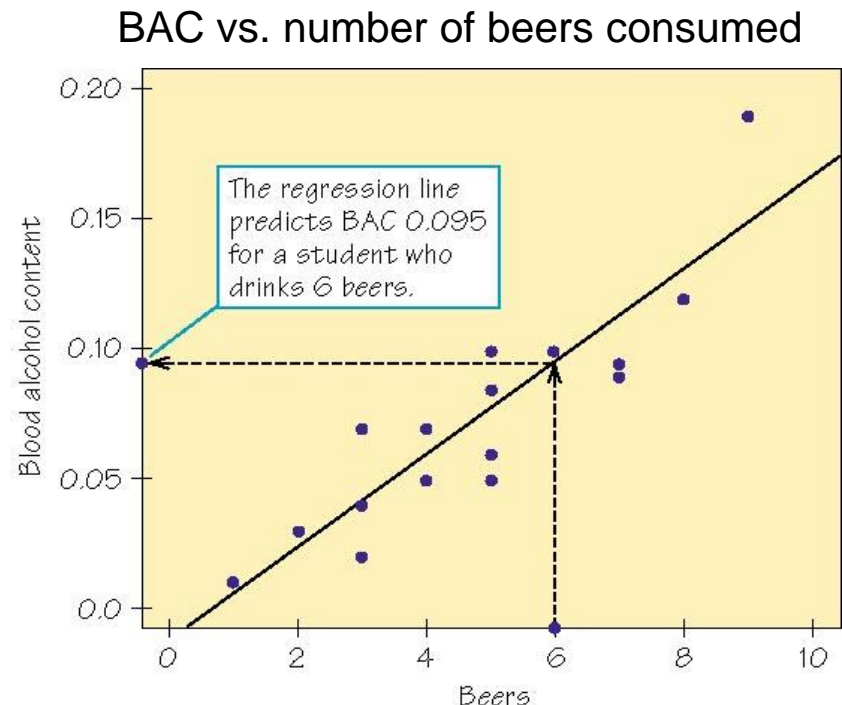
■ Regression Line

- A straight line that describes how a response variable y changes as an explanatory variable x changes.
- Regression lines are often used to predict the value of y for a given value of x .

A regression line has been added to be able to predict blood alcohol content from the number of beers a student drinks.

Graphically, you can predict that if $x = 6$ beers, then $y = 0.095$ BAC.

(Legal limit for driving in many states is $BAC = 0.08$.)



Chapter 6: Exploring Data: Distributions

Regression Lines

■ Using the Equation of the Line for Predictions

- It is easier to use the equation of the line for predicting the value of y , given the value of x .

Using the equation of the line for the previous example:

$$\text{predicted BAC} = -0.0127 + (0.01796)(\text{beers})$$

$$y = -0.127 + 0.01796 (x)$$

For a young adult drinking 6 beers ($x = 6$):

$$\text{predicted BAC} = -0.0127 + 0.01796 (6) = 0.095$$

■ Straight Lines

- A straight line for predicting y from x has an equation of the form:

$$\hat{y} = mx + b$$

- In this equation, m is the slope, the amount by which y changes when x increases by 1 unit.
- The number b is the y -intercept, the value of y when $x = 0$.

Chapter 6: Exploring Data: Distributions

Correlation

■ Correlation, r

- Measures the direction and strength of the straight-line relationship between two numerical variables.
- A correlation r is always a number between -1 and 1 .
- It has the same sign as the slope of a regression line.
 - $r > 0$ for positive association (increase in one variable causes an increase in the other).
 - $r < 0$ for negative association (increase in one variable causes a decrease in the other)

Chapter 6: Exploring Data: Distributions

Correlation

■ Correlation, r

- ❑ Perfect correlation $r = 1$ or $r = -1$ occurs only when all points lie exactly on a straight line.
- ❑ The correlation moves away from 1 or -1 (toward zero) as the straight-line relationship gets weaker.
- ❑ Correlation $r = 0$ indicates no straight-line relationship.

Chapter 6: Exploring Data: Distributions

Correlation

■ Correlation

- Correlation is strongly affected by a few outlying observations. (Also, the mean and standard deviation are affected by outliers.)

■ Equation of the Correlation

- To calculate the correlation, suppose you have data on variable \mathbf{x} and \mathbf{y} for n individuals.
- From the data, you have the values calculated for the means and standard deviations for \mathbf{x} and \mathbf{y} .
 - The means and standard deviations for the two variables are \bar{x} and s_x for the \mathbf{x} -values, and \bar{y} and s_y for the \mathbf{y} -values.
- The correlation r between \mathbf{x} and \mathbf{y} is:

$$r = \frac{1}{n-1} \left[\frac{(x_1 - \bar{x})(y_1 - \bar{y})}{s_x s_y} + \frac{(x_2 - \bar{x})(y_2 - \bar{y})}{s_x s_y} + \dots + \frac{(x_n - \bar{x})(y_n - \bar{y})}{s_x s_y} \right]$$

Chapter 6: Exploring Data: Distributions

Correlation

■ Correlation

- The scatterplots below show examples of how the correlation r measures the direction and strength of a straight-line association.



Correlation $r = 0$



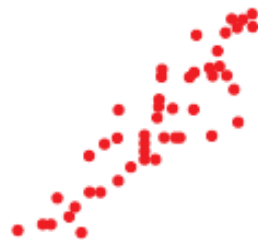
Correlation $r = -0.3$



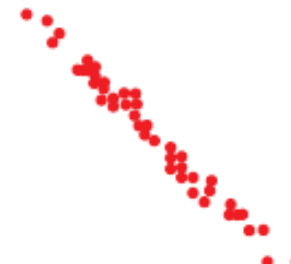
Correlation $r = 0.5$



Correlation $r = -0.7$



Correlation $r = 0.9$



Correlation $r = -0.99$

Chapter 6: Exploring Data: Distributions

Least-Squares Regression

■ Least-Squares Regression Line

- A line that makes the sum of the squares of the vertical distances of the data points from the line as small as possible.

■ Equation of the Least-Squares Regression Line

- From the data for an explanatory variable x and a response variable y for n individuals, we have calculated the means \bar{x} , \bar{y} , and standard deviations s_x , s_y , as well as their correlation r .

The least-squares regression line is the line:

$$\text{Predicted } \hat{y} = mx + b$$

$$\text{With slope ... } m = r \frac{s_y}{s_x}$$

$$\text{And y-intercept ... } b = \bar{y} - m\bar{x}$$

This equation was used to calculate the line for predicting BAC for number of beers drunk.

$$\text{Predicted } y = -0.0127 + 0.01796 x$$

Chapter 6: Exploring Data: Distributions

Interpreting Correlation and Regression

- A Few Cautions When Using Correlation and Regression
 - Both the correlation r and least-squares regression line can be strongly influenced by a few outlying points.
 - Always make a scatterplot before doing any calculations.
 - Often the relationship between two variables is strongly influenced by other variables.
 - Before conclusions are drawn based on correlation and regression, other possible effects of other variables should be considered.

Chapter 6: Exploring Data: Distributions

Interpreting Correlation and Regression

- A Few Cautions When Using Correlation and Regression
 - A strong association between two variables is not enough to draw conclusions about cause and effect.
 - Sometimes an observed association really does reflect cause and effect (such as drinking beer causes increased BAC).
 - Sometimes a strong association is explained by other variables that influence both x and y .
 - *Remember, association does not imply causation.*