Chapter 10: Correlation and Regression

Scatterplot--a two-dimensional graph of data values

- Use a *scatterplot* to look at the relationship between two quantitative variables
- Plot has one variable's values along the vertical axis and the other variable's values along the horizontal axis
- Simple Linear Regression—a method to generate linear equation that describes the average relationship between a response and explanatory variable
- **Correlation**, a statistic that measures the *strength* and *direction* of a linear relationship



Relationships Between Quantitative Variables

Questions that might be asked

- 1. What is the *average* pattern? Does it look like a straight line or is it curved?
- 2. What is the direction of the pattern?
- 3. How much do individual points vary from the average pattern?

3

4. Are there any unusual data points?

Relationships Between Quantitative Variables

Positive/Negative Association

- Positive association--the values of one variable tend to increase as the values of the other variable increase.
- Negative association--the values of one variable tend to decrease as the values of the other variable increase.















Correlation Coefficient

Correlation is a measure of the linear association between two variables, *x* and *y*.

Pearson product moment coefficient of correlation, r, is a measure of the strength of the linear relationship between two variables x and y. It is computed for a **sample** of n measurements on x and y as:

$$r = \frac{1}{n-1} \sum_{i} \left(\frac{x_i - \overline{x}}{s_x} \right) \left(\frac{y_i - \overline{y}}{s_y} \right)$$





Correlation Coefficient

- 1. *r* is positive when the slope is positive and likewise negative when the slope is negative
- 2. If $SS_{xy} = 0$, then $r = 0 \Rightarrow$ that there is no association between the magnitudes of the two variables; or, a change in the magnitude of one variable does not imply a change in the magnitude of the other variable
- 3. The correlation coefficient, *r*, is unitless and assumes a value between -1 and +1, regardless of the units of *x* and *y*

13

4. The Cor(X, Y) = Cor(Y, X)



Wing (x)	Tail (y)	x^2	y^2	xy
10.4	7.4	108.16	54.76	76.96
10.8	7.6	116.64	57.76	82.08
11.1	7.9	123.21	62.41	87.69
10.2	7.2	104.04	51.84	73.44
10.3	7.4	106.09	54.76	76.22
10.2	7.1	104.04	50.41	72.42
10.7	7.4	114.49	54.76	79.18
10.5	7.2	110.25	51.84	75.6
84.2	59.2	886.92	438.54	623.59









Correlation Coefficient--Interpretations

Interpretations of an observed association (continued)

- 3. There is no causation. The association is explained by how the explanatory and response variables are both affected by other variables.
- 4. The response variable is causing a change in the explanatory variable.

Correlation Coefficient

Population Correlation Coefficient

Denoted by $\rho(rho)$

 ρ is estimated by the sample statistic, r

Hypothesis test for ρ given on page 527—we will not specifically go over this

19

21

Simple Linear Regression

Regression equation—an equation that describes the average relationship between a response (dependent) and an explanatory (independent) variable.



Deterministic Model

A model that defines an exact relationship between variables.

Example: y = 1.5x

There is no allowance for error in the prediction of y for a given x.





20

Probabilistic Model

A model that accounts for *random error*.

Includes both a deterministic component and a random error component.

y = 1.5x + random error

This model hypothesizes a probabilistic relationship between *y* and *x*.



Probabilistic Model—General Form y = Deterministic component + Random component where y is the "variable of interest". Assume that the mean value of the random error is zero → the mean value of y, E(y), equals the deterministic component of the model

$y = \beta_0 + \beta_1 x + \varepsilon$	where <i>y</i> = <i>Dependent variable</i> <i>x</i> = <i>Independent variable</i>
$eta_0 = population y$ -int	ercept of the line—the point at which the line intersects or cuts through the y-axis
$\beta_1 = population slope$	of the line—the amount of increase (or decrease) in the deterministic component of y for every 1-unit increase (or decrease) in x.
\mathcal{E} = random error com	ponent

First-Order (Straight Line) Probabilistic Model

 β_0 and β_1 are population parameters. They will only be known if the population of all (x, y) measurements are available.

 β_0 and β_1 , along with a specific value of the independent variable *x* determine the *mean value* of the dependent variable *y*.

29

27

25

Model Development

 β_0 and β_1 will generally be unknown.

The process of developing a model, estimating model parameters, and using the model can be summarized in these 5-steps:

1. Hypothesize the deterministic component of the model that relates the mean, *E*(*y*) to the independent variable *x*

$$E(y) = \beta_0 + \beta_1 x$$

2. Use sample data to estimate unknown model parameters

find estimates: $\hat{\beta}_0$ or b_0 , $\hat{\beta}_1$ or b_1



Model Development (continued)

3. Specify the probability distribution of the random error term and estimate the SD of this distribution

 $\varepsilon_i \sim N(0,\sigma)$ – -will revisit this later

- 4. Statistically evaluate the usefulness of the model
- 5. Use model for prediction, estimation or other purposes

31

Subject	Amount of Drug (%) x	Reaction Time (seconds) V
1	1	1
2	2	1
3	3	2
4	4	2
5	5	4





Errors the ob	s <i>of p</i> serve	o <i>red</i> ed a	<i>liction</i> ver nd the predi	tical different cted values	nces between of <i>y</i>	
	x	y	$\tilde{y} = -1 + x$	$(y-\tilde{y})$	$(y-\tilde{y})^2$	
	1	1	0	(1-0) = 1	1	
	2	1	1	(1-1) = 0	0	
	3	2	2	(2-2) = 0	0	
	4	2	3	(2-3) = -1	1	
	5	4	4	(4-4) = 0	0	
				Sum of	Sum of squared	
				errors = 0	errors (SSE) = 2	



 b_0 and b_1 are estimators of β_0 and β_1



Least Squares Line (continued)

Define the sum of squares of the deviations of the y values about their predicted values for all n data points as:

SSE =
$$\sum_{i=1}^{n} (y_i - \hat{y})^2 = \sum_{i=1}^{n} [y_i - (b_0 + b_1 x_i)]^2$$

We want to find b_0 and b_1 to make the SSE a minimum---termed *least squares estimates*

 $\hat{y} = b_0 + b_1 x$ is called the least squares line

Formulas for the Least Squares Estimates

$$Slope: b_{1} = \frac{SS_{xy}}{SS_{xx}} \text{ or } b_{1} = r \frac{SD_{y}}{SD_{x}}$$

$$SS_{xy} = \sum (x_{i} - \bar{x})(y_{i} - \bar{y}) \qquad SS_{xx} = \sum (x_{i} - \bar{x})^{2}$$

$$= \sum x_{i}y_{i} - \frac{(\sum x_{i})(\sum y_{i})}{n} \qquad = \sum x_{i}^{2} - \frac{(\sum x_{i})^{2}}{n}$$

$$y\text{-intercept: } b_{0} = \overline{y} - b_{1}\overline{x} = \frac{\sum y_{i}}{n} - b_{1}\frac{\sum x_{i}}{n}$$

$$n = \text{ sample size}$$

$$38$$





x	у	$\hat{y} =1 + .7x$	$(y-\hat{y})$	$(y-\hat{y})^2$
1	1	.6	(16) = .4	.16
2	1	1.3	(1-1.3) =3	.09
3	2	2.0	(2-2.0) = 0	.00
4	2	2.7	(2-2.7) =7	.49
5	4	3.4	(4-3.4) = .6	.36
			Sum of errors = 0	Sum of squared errors (SSE) = 1.10

41





Least Squares Line—Interpretation of $\hat{y} = -.1 + .7x$

The slope of 0.7 implies that for every unit increase of x, the *mean value* of y is estimated to increase by 0.7 units.

In the context of the problem:

For every 1% increase in the amount of drug in the bloodstream, the mean reaction time is estimated to increase by 0.7 seconds over the sampled range of drug amounts from 1% to 5%.





Coefficient of Determination (continued)	
$SS_{yy} = \sum (y_i - \overline{y})^2$ total sample variation around mean	
$SSE = \sum (y_i - \hat{y}_i)^2unexplained sample variability after fitting$	
SS_{yy} – SSE explained sample variability attributable to linear relationship	
$\frac{SS_{yy} - SSE}{SS_{yy}} = \frac{explained}{total} = \frac{explained}{variability} = \frac{explained}{variability}$	
	46



