

# Chapter 10: Correlation and Regression

## Section 10.1: Correlation

Stat 50

### CORRELATION

*Def* A correlation exists between two variables when the values of one variable are somehow associated with the values of the other variable.

### EXPLANATORY VS RESPONSE VARIABLE

*Def* Explanatory Variable

*Def* Response Variable

The way to distinguish the difference between these two variables is by asking, “Which statement makes sense?”

EX: A researcher wants to examine whether babies fed on breast milk are more or less likely to be ill.

(a) Feeding a baby on breast milk causes resistance to disease.

(b) Resistance to disease causes a baby to feed on breast milk.

Can you identify which variable is which? – presence or absence of breast milk – resistance to disease

EX: Identify the explanatory and response variables in the following examples.

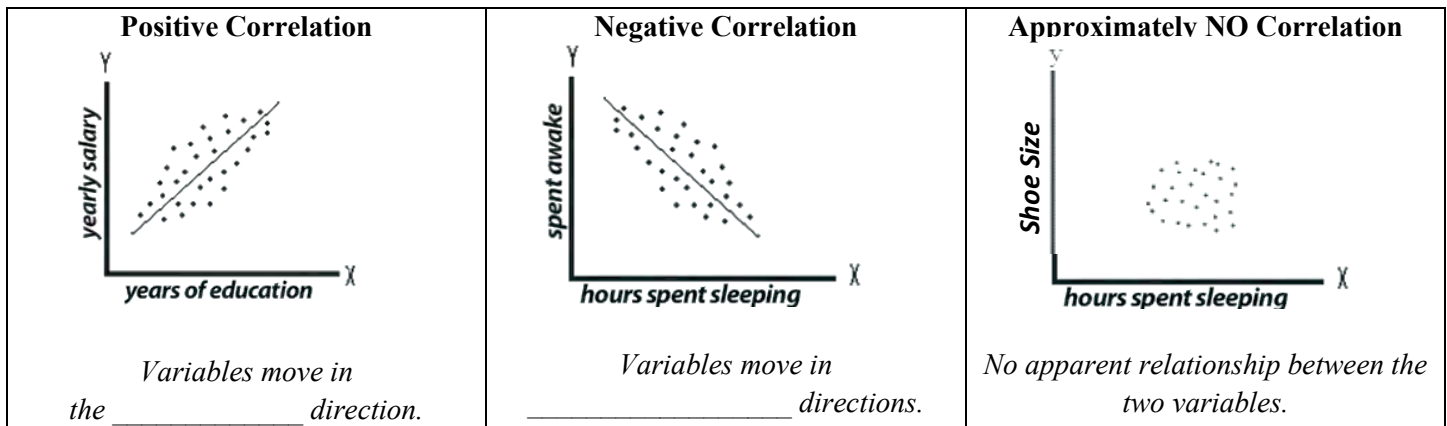
(a) An experiment was conducted to test the effects of sleep deprivation on human response times.

(b) Researcher Penny Gordon Larson and her associate wanted to determine whether young couples who marry or cohabitate are more likely to gain weight than those who stay single.

### SCATTERPLOTS

*Def* A **scatterplot** is a plot of paired  $(x, y)$  quantitative data.

*Note:* A scatter diagram is often helpful in determining whether there is a relationship between the two variables.



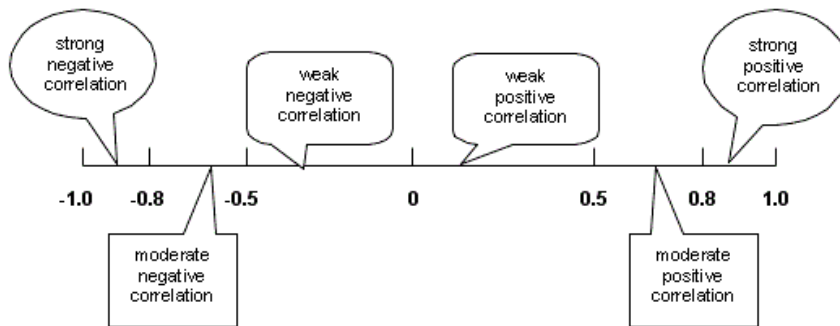
To test for \_\_\_\_\_ correlation we attempt to draw a line that best fits the data. Every linear correlation is expressed by two features:

STRENGTH –	DIRECTION –
------------	-------------

The strength of the linear correlation is represented by a numerical value called the \_\_\_\_\_.

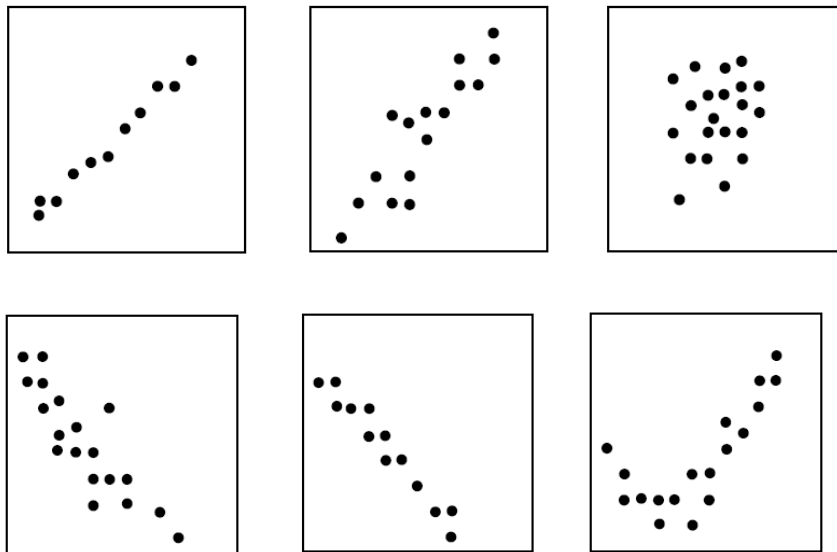
**PROPERTIES OF THE LINEAR CORRELATION COEFFICIENT**

1. The value of  $r$  is always between  $-1$  and  $1$ , inclusive.  $-1 \leq r \leq 1$
2. The value of  $r$  does not change if all values of either variable are converted to a different scale.
3.  $r$  measures the strength of a linear relationship only. (It does not measure nonlinear relationships.)



\*Note: Correlation does not imply \_\_\_\_\_, just \_\_\_\_\_.

EX: Choose from the following word bank to determine the features of each scatter plot below.



Choose a word from each category to describe each scatter plot shown to the left and write near/next to each plot:

---

- Linear relationship
- Non-Linear Relationship
- No Relationship

---

- Perfect
- Strong
- Moderate
- Weak

---

- Positive Association
- Negative Association
- No Correlation

---

- $r = -0.735$
- $r = 0.634$
- $r = -0.02$
- $r = 0.886$
- $r = -0.89$
- *no r value*

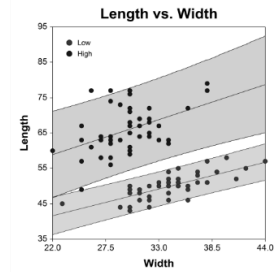
**EXPLAINED VARIATION**

The value of  $r^2$  is the proportion of the variation in  $y$  that can be explained by the linear relationship between  $x$  and  $y$ .

Different samples will produce different scatterplots, and thus, different  $r$  values.

Our job is to take a large enough sample to get close to the *true* population linear correlation coefficient called \_\_\_\_\_.

We are going to assume there is no correlation ( $\rho = \text{_____}$ ) and try to prove otherwise based on our sample.



Steps for Hypothesis Test when Applied to testing $\rho$		
<p><b>Check Requirements</b></p> <ul style="list-style-type: none"> <li>• Simple Random Sample</li> <li>• Visual examination shows straight line pattern</li> <li>• Remove any outliers</li> </ul>	<p><b>Step 1: Hypotheses</b></p> <p style="text-align: center;"><math>H_0: \rho = 0</math> (there is no linear correlation)</p> <p style="text-align: center;"><math>H_1: \rho \neq 0</math> (there <b>is</b> a linear correlation)</p>	<p><b>Step 2: Level of Significance</b></p>
<p><b>Step 3: Test Statistic</b></p> <div style="text-align: center; margin: 10px 0;"> <math display="block">t_0 = \frac{r}{\sqrt{\frac{1-r^2}{df}}}</math> </div> <p style="text-align: right; margin-right: 100px;">where <math>df = n - 2</math></p>		
<p><b>Step 4: Find a Critical Value or P-Value</b></p>		
<p><b>P-VALUE METHOD</b></p>	<p><b>DECISION</b></p>	<p style="font-size: 1.2em;">{</p> <p>Reject <math>H_0 \sim</math> if <math>P\text{-value} \leq \alpha</math></p> <p>Fail to Reject <math>H_0 \sim</math> if <math>P\text{-value} &gt; \alpha</math></p>
<p><b>CRITICAL VALUE METHOD</b></p>	<p><b>DECISION</b></p>	<p style="font-size: 1.2em;">{</p> <p>Reject <math>H_0 \sim</math> if <math>r^*</math> lies in the critical region</p> <p>Fail to Reject <math>H_0 \sim</math> if <math>r^*</math> doesn't lie in the critical region</p>
<p><b>Step 5: Write a CONCLUSION</b> either rejecting or failing to reject <math>H_0</math></p>		

$n$	$\alpha = .05$	$\alpha = .01$
4	.950	.990
5	.878	.959
6	.811	.917
7	.754	.875
8	.707	.834
9	.666	.798
10	.632	.765
11	.602	.735
12	.576	.708
13	.553	.684
14	.532	.661
15	.514	.641

$n$	$\alpha = .05$	$\alpha = .01$
16	.497	.623
17	.482	.606
18	.468	.590
19	.456	.575
20	.444	.561
25	.396	.505
30	.361	.463
35	.335	.430
40	.312	.402
45	.294	.378
50	.279	.361
60	.254	.330

$n$	$\alpha = .05$	$\alpha = .01$
70	.236	.305
80	.220	.286
90	.207	.269
100	.196	.256

**Critical Values of the Pearson Correlation Coefficient  $r$**

GRAPHING CALCULATOR  
(TI-83 OR 84)

Instructions:

STAT  $\Rightarrow$  TESTS  $\Rightarrow$  LinRegTTest

Ex: The following table gives information on average saturated fat (in grams) consumed per day and cholesterol level (in milligrams per centiliters) of ten men taken from a simple random sample.

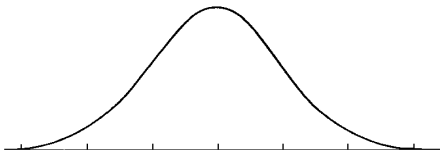
	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>10</i>
Fat Consumption (in grams)	55	68	50	34	43	58	77	36	60	39
Cholesterol level (in mg/cL)	180	215	195	165	170	204	235	150	190	185

Use a 0.01 significance level to determine if there is a linear correlation between saturated fat consumption and cholesterol level.

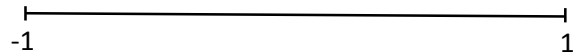
*Null and Alternative Hypothesis*

*Test Statistic*

*P-value:*



*Critical Value:*



*Decision about Null Hypothesis*

*Conclusion*

Ex: The following table gives the total 2004 payroll (on the opening day of the season, rounded to the nearest million dollars) and the percentage of games won in 2004 by each National League team.

	<i>D'backs</i>	<i>Braves</i>	<i>Cubs</i>	<i>Reds</i>	<i>Rockies</i>	<i>Marlins</i>	<i>Astros</i>	<i>Dodgers</i>
Payroll (in millions)	70	90	91	47	65	42	75	93
Percentage of Wins	31.5	59.3	54.9	46.9	42.0	51.2	56.8	57.4

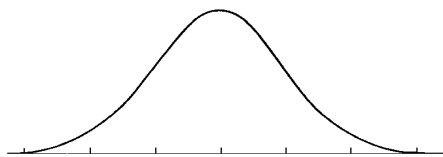
	<i>Brewers</i>	<i>Expos</i>	<i>Mets</i>	<i>Phillies</i>	<i>Pirates</i>	<i>Cards</i>	<i>Padres</i>	<i>Giants</i>
Payroll (in millions)	28	41	97	93	32	83	55	82
Percentage of Wins	41.6	41.4	43.8	53.1	44.7	64.8	53.7	56.2

Use a 0.05 significance level to determine if there's a correlation between payroll and percentage of games won.

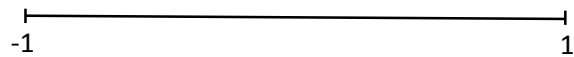
*Null and Alternative Hypothesis*

*Test Statistic*

*P-value:*



*Critical Value:*



*Decision about Null Hypothesis*

*Conclusion*

**REGRESSION**

*Def* Given a collection of paired sample data, the **regression equation**  $\hat{y} = b_0 + b_1x$  algebraically describes the relationship between the two variables.

*Note:* The graph of the regression equation is called the **regression line** or *line of best-fit*

**TERMINOLOGY**

$x$  is referred to as the explanatory variable, predictor variable, or the independent variable.

$y$  is referred to as the response variable or the dependent variable.

<b>REGRESSION LINE</b> <b>A.K.A.</b> <b>LEAST-SQUARES LINE</b> <b>A.K.A.</b> <b>"LINE OF BEST FIT"</b>	The line that minimizes the distance between the points and the line. It _____ _____ the data points. $\hat{y} = b_0 + b_1x$ where the point $(\bar{x}, \bar{y})$ will always be on the least-squares line. $b_1 =$ _____ $b_0 =$ _____
--	--

**FORMULAS**

**Slope**  $b_1 = r \cdot \frac{s_y}{s_x}$       **y-intercept**  $b_0 = \bar{y} - b_1\bar{x}$

*Round-Off Rule:* Round the slope and y-intercept to three significant digits.

EX: Below is a sample of five patients at a hospital with the information regarding their height and weight.

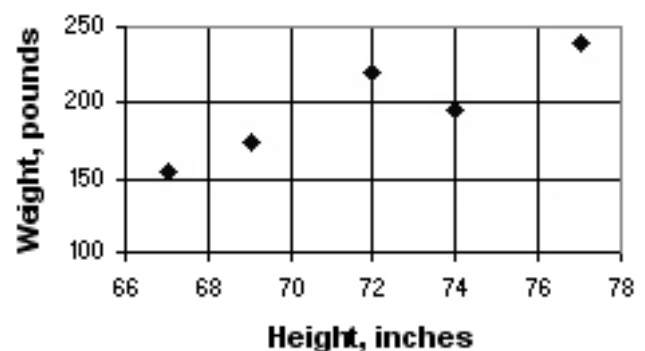
(a) Describe the relationship of the data using the scatterplot given.

(b) Find the correlation coefficient/test statistic and determine whether a correlation exists.

(c) Find the line of best fit and sketch it below.

(c) Interpret the slope.

Height (in)	Weight (lbs)
67	155
72	220
77	240
74	195
69	175



**Interpretation of slope:** For every unit increase in the explanatory variable, on average, there is an increase/decrease of “slope” units on the response variable.

**Interpretation of y-intercept:** When the explanatory variable is 0 units, on average, the response variable is \_\_\_\_\_ units.

GRAPHING CALCULATOR (TI-83 OR 84)

Instructions: STAT  $\Rightarrow$  TESTS  $\Rightarrow$  LinRegTTest

EX: Recall the exercise from the last section in which we concluded there was a significant linear correlation between the average saturated fat consumed per day and the cholesterol level of ten men.

	1	2	3	4	5	6	7	8	9	10
Fat Consumption (in grams)	55	68	50	34	43	58	77	36	60	39
Cholesterol level (in mg/cL)	180	215	195	165	170	204	235	150	190	185

- (a) Find the regression equation where  $x$  is the average daily fat consumption (in grams) of a man and  $y$  is the cholesterol level (in mg/cL).
- (b) Interpret the slope in the context of the problem.
- (c) Predict the cholesterol level of a man who consumes 65 grams of saturated fat per day.

### **PREDICTIONS**

If there is a significant linear correlation between  $x$  and  $y$ , then use the \_\_\_\_\_ to predict the value of  $y$  given a specific value of  $x$ .

If there is no significant linear correlation between  $x$  and  $y$ , then the best prediction of  $y$  is the \_\_\_\_\_ for any given value of  $x$ .

EX: Are fat and sodium content related in fast food? Here are the fat and sodium content for several brands of burgers.

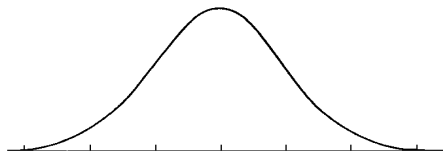
	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>
Fat (in grams)	19	31	34	35	39	39	43
Sodium(mg)	920	1500	1310	860	1180	940	1260

Use a 0.05 significance level to determine if there is a linear correlation between fat and sodium content in burgers.

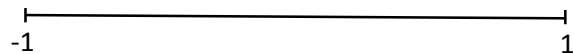
*Null and Alternative Hypothesis*

*Test Statistic (or correlation coefficient)*

*P-value:*



*Critical Value:*



*Decision*

*Conclusion*

(a) What is the regression equation? Is it helpful in this situation? Why or why not?

(a) Predict the sodium level of a burger with 25 grams of fat.

GRAPHING CALCULATOR (TI-83 OR 84)

To create and view a Scatterplot and Linear Regression Line

Instructions:

- 1) 2<sup>nd</sup> ⇒ 0 (catalog) ⇒ DiagnosticOn ⇒ Enter
- 2) STAT ⇒ EDIT (enter 1<sup>st</sup> Variable in L<sub>1</sub> and 2<sup>nd</sup> Variable in L<sub>2</sub>)
- 3) STAT ⇒ CALC ⇒ 4: LinReg ( $ax + b$ ) ⇒ Store RegEQ: ⇒ Vars ⇒ Y-Vars ⇒ 1: Function ⇒ 1:Y<sub>1</sub> ⇒ Calculate
- 4) 2<sup>nd</sup> ⇒ y = ⇒ 1: Plot1 ⇒ On ⇒ Zoom ⇒ 9: ZoomStat