

# Stat 13, Intro. to Statistical Methods for the Life and Health Sciences.

1. Two variables and correlation.
2. Simulation based inference for correlation.
3. Calculating correlation,  $r$ .
4. Linear regression.

HW3, due Fri, May22, 1159pm. 4.CE.10, 5.3.28, 6.1.17, and 6.3.14.

In 5.3.28d, use the theory-based formula. You do not need to use an applet.

Read ch7 and 10.

The course website is <http://www.stat.ucla.edu/~frederic/13/S26> .

# 1. Two quantitative variables and correlation.

Chapter 10

# Two Quantitative Variables: Scatterplots and Correlation

Section 10.1

# Scatterplots and Correlation

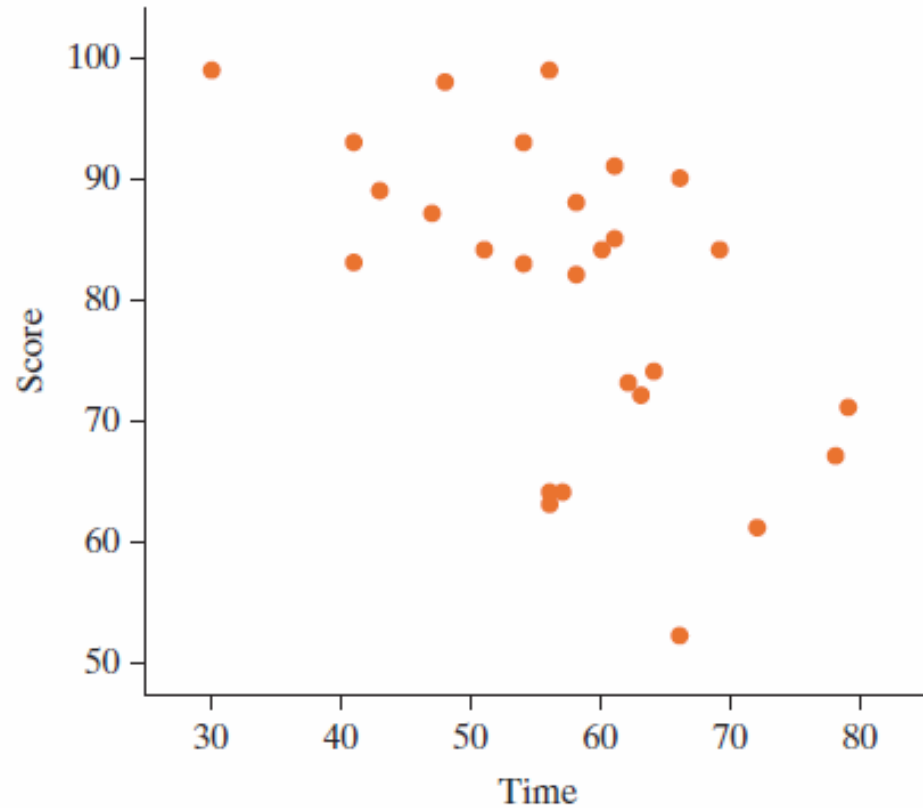
Suppose we collected data on the relationship between the time it takes a student to take a test and the resulting score.

<b>Time</b>	30	41	41	43	47	48	51	54	54	56	56	56	57	58
<b>Score</b>	100	84	94	90	88	99	85	84	94	100	65	64	65	89
<b>Time</b>	58	60	61	61	62	63	64	66	66	69	72	78	79	
<b>Score</b>	83	85	86	92	74	73	75	53	91	85	62	68	72	

# Scatterplot

Put explanatory variable on the horizontal axis.

Put response variable on the vertical axis.



# Describing Scatterplots

- When we describe data in a scatterplot, we describe the
  - Direction (positive or negative)
  - Form (linear or not)
  - Strength (strong-moderate-weak, we will let correlation help us decide)
  - Unusual Observations
- How would you describe the time and test scatterplot?

# Correlation

- **Correlation** measures the strength and direction of a linear association between two quantitative variables.
- Correlation is a number between -1 and 1.
- With positive correlation one variable increases, on average, as the other increases.
- With negative correlation one variable decreases, on average, as the other increases.
- The closer it is to either -1 or 1 the closer the points fit to a line.
- The correlation for the test data is -0.56.

# Correlation Guidelines

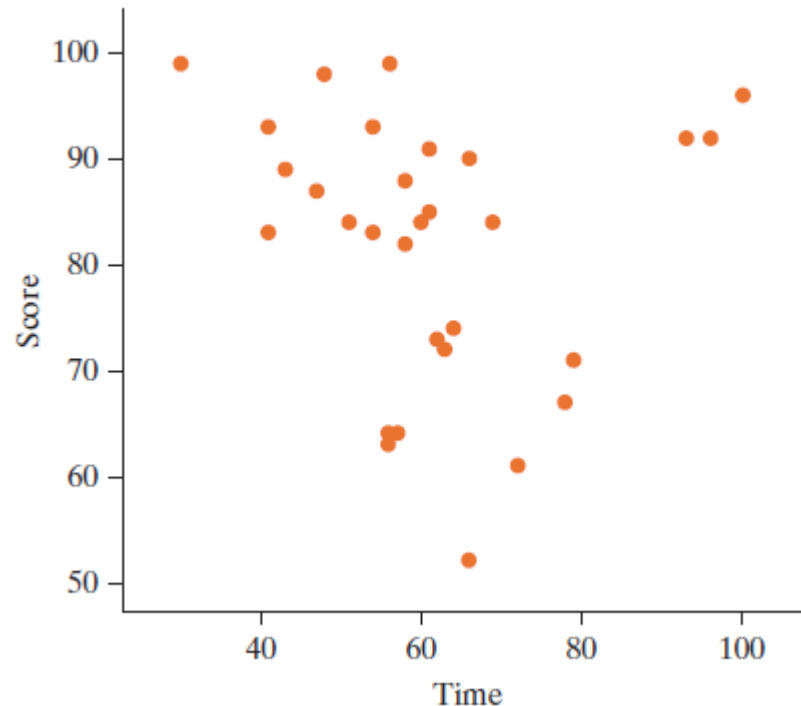
Correlation Value	Strength of Association	What this means
0.7 to 1.0	Strong	The points will appear to be nearly a straight line
0.3 to 0.7	Moderate	When looking at the graph the increasing/decreasing pattern will be clear, but there is considerable scatter.
0.1 to 0.3	Weak	With some effort you will be able to see a slightly increasing/decreasing pattern
0 to 0.1	None	No discernible increasing/decreasing pattern

Same Strength Results with Negative Correlations

# Back to the test data

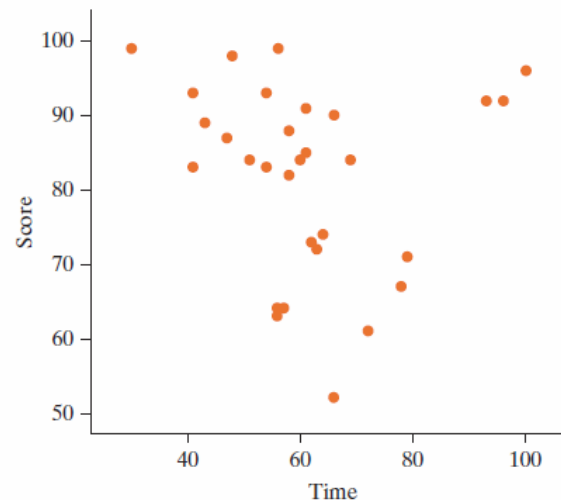
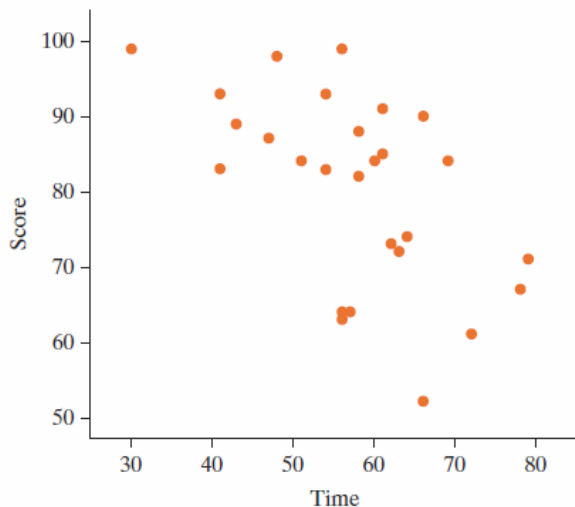
Actually the last three people to finish the test had scores of 93, 93, and 97.

What does this do  
to the correlation?



# Influential Observations

- The correlation changed from  $-0.56$  (a fairly moderate negative correlation) to  $-0.12$  (a weak negative correlation).
- Points that are far to the left or right and not in the overall direction of the scatterplot can greatly change the correlation. (influential observations)



# Correlation

- **Correlation** measures the strength and direction of a linear association between two quantitative variables.
  - $-1 \leq r \leq 1$
  - Correlation makes no distinction between explanatory and response variables.
  - Correlation has no units.
  - Correlation is not resistant to outliers. It is sensitive.

# Learning Objectives for Section 10.1

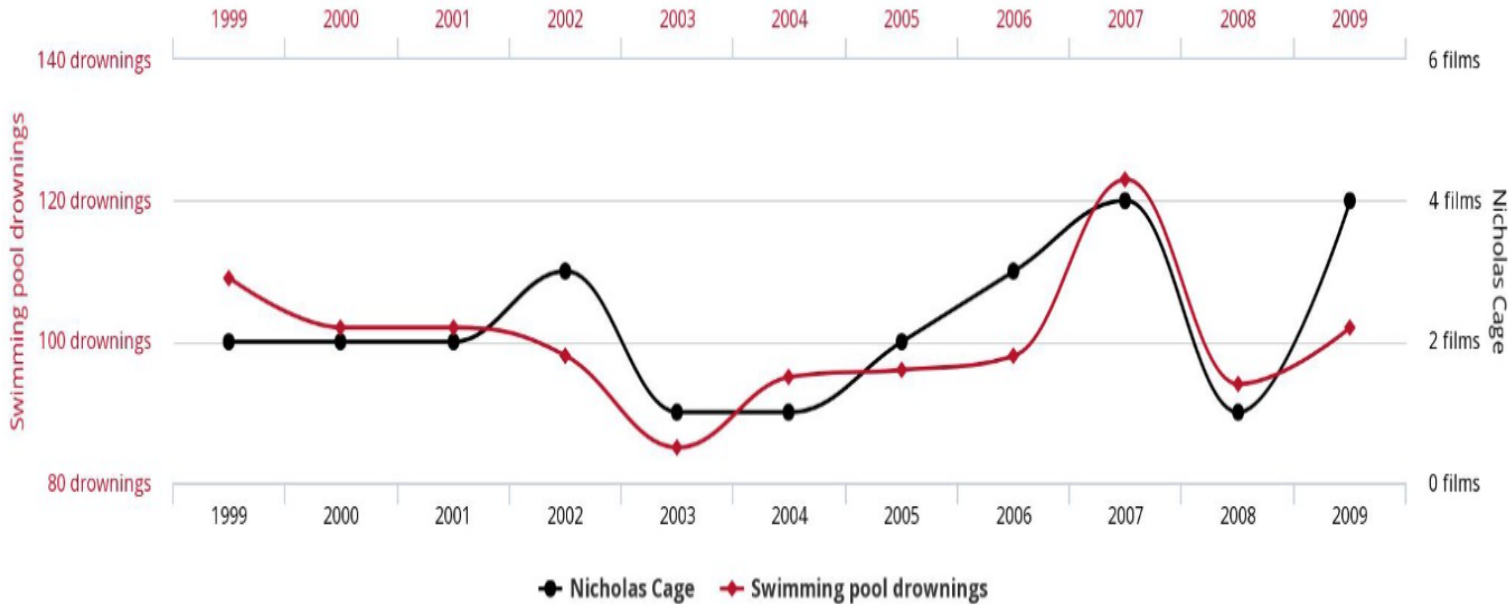
- Summarize the characteristics of a scatterplot by describing its direction, form, strength and whether there are any unusual observations.
- Recognize that the correlation coefficient is appropriate only for summarizing the strength and direction of a scatterplot that has linear form.
- Recognize that a scatterplot is the appropriate graph for displaying the relationship between two quantitative variables and create a scatterplot from raw data.
- Recognize that a correlation coefficient of 0 means there is no linear association between the two variables and that a correlation coefficient of -1 or 1 means that the scatterplot is exactly a straight line.
- Understand that the correlation coefficient is influenced by extreme observations.

# Note that correlation $\neq$ causation.

Number of people who drowned by falling into a pool

correlates with

Films Nicolas Cage appeared in



tylervigen.com

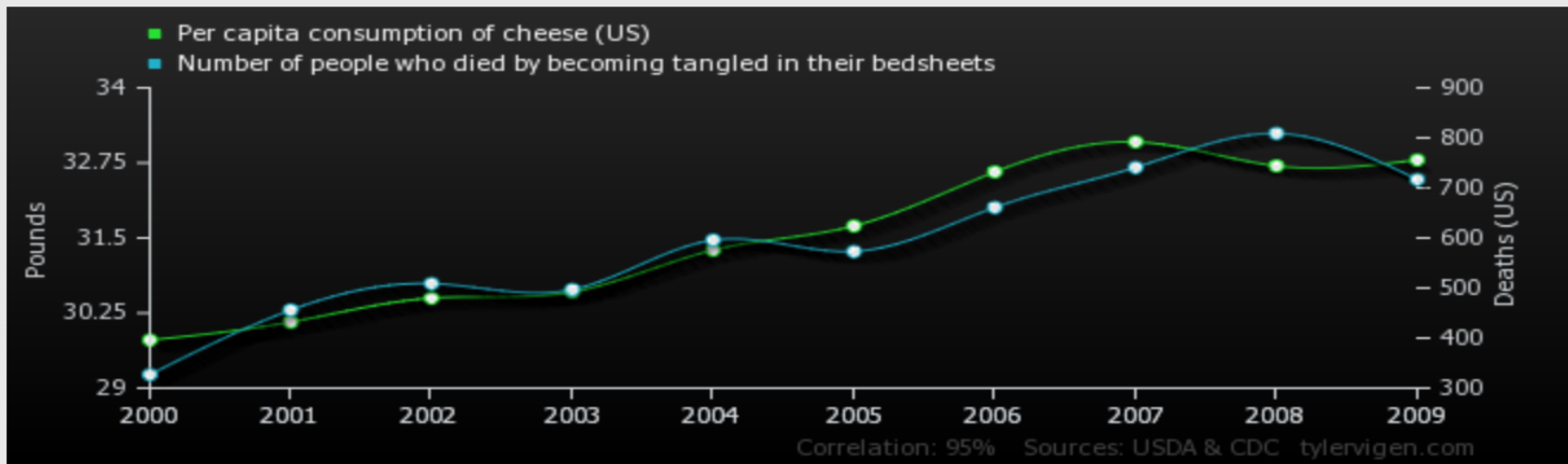
from: <http://tylervigen.com>

# Note that correlation $\neq$ causation.

## Per capita consumption of cheese (US)

correlates with

## Number of people who died by becoming tangled in their bedsheets



	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009
<i>Per capita consumption of cheese (US) Pounds (USDA)</i>	29.8	30.1	30.5	30.6	31.3	31.7	32.6	33.1	32.7	32.8
<i>Number of people who died by becoming tangled in their bedsheets Deaths (US) (CDC)</i>	327	456	509	497	596	573	661	741	809	717

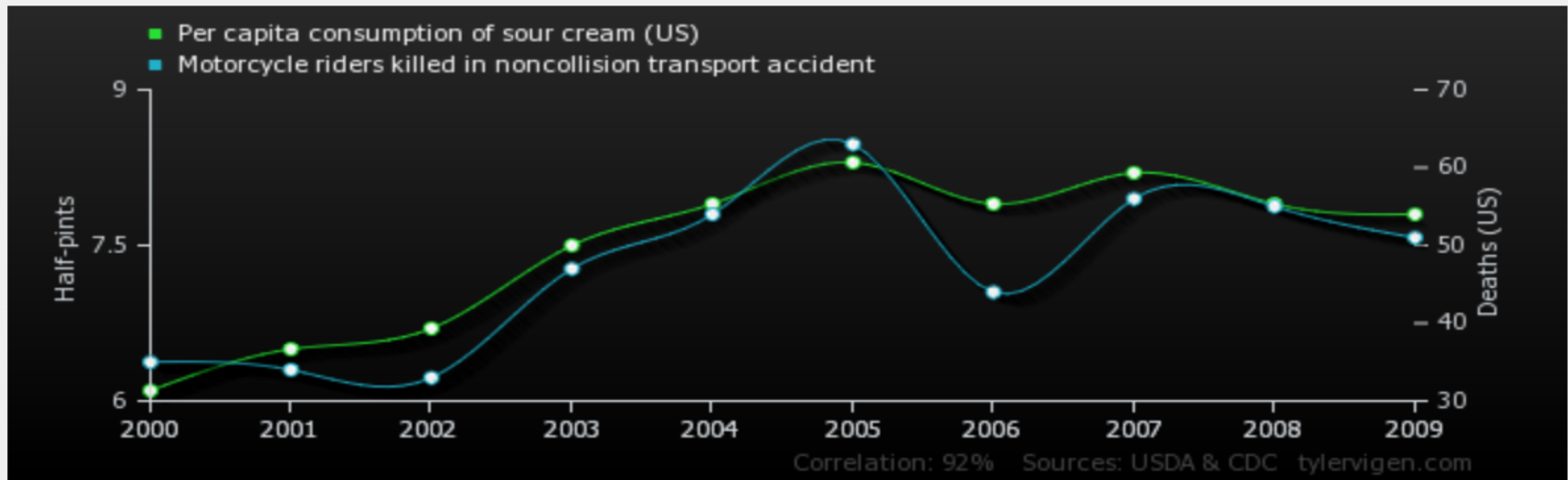
**Correlation: 0.947091**

# Note that correlation $\neq$ causation.

## Per capita consumption of sour cream (US)

correlates with

## Motorcycle riders killed in noncollision transport accident



	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009
<i>Per capita consumption of sour cream (US) Half-pints (USDA)</i>	6.1	6.5	6.7	7.5	7.9	8.3	7.9	8.2	7.9	7.8
<i>Motorcycle riders killed in noncollision transport accident Deaths (US) (CDC)</i>	35	34	33	47	54	63	44	56	55	51

**Correlation: 0.916391**

# Inference for the Correlation Coefficient: Simulation-Based Approach

Section 10.2

We will look at a small sample example to see if body temperature is associated with heart rate.

# Temperature and Heart Rate

## Hypotheses

- Null: There is no association between heart rate and body temperature. ( $\rho = 0$ )
- Alternative: There is a positive linear association between heart rate and body temperature. ( $\rho > 0$ )

$\rho = \text{rho}$

# Inference for Correlation with Simulation

(Section 10.2)

1. Compute the observed statistic. (Correlation)
2. Scramble the response variable, compute the simulated statistic, and repeat this process many times.
3. Reject the null hypothesis if the observed statistic is in the tail of the null distribution.

# Temperature and Heart Rate

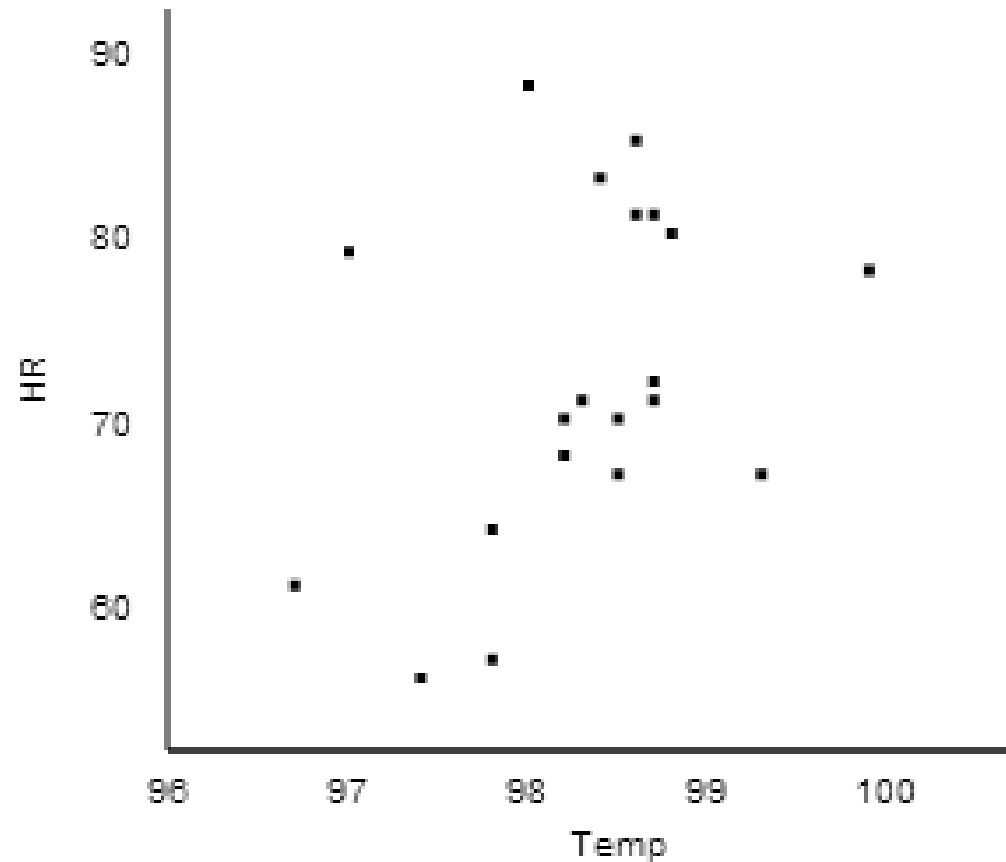
Collect the Data

<b>Tmp</b>	98.3	98.2	98.7	98.5	97.0	98.8	98.5	98.7	99.3	97.8
<b>HR</b>	72	69	72	71	80	81	68	82	68	65
<b>Tmp</b>	98.2	99.9	98.6	98.6	97.8	98.4	98.7	97.4	96.7	98.0
<b>HR</b>	71	79	86	82	58	84	73	57	62	89

# Temperature and Heart Rate

Explore the Data

$$r = 0.378$$



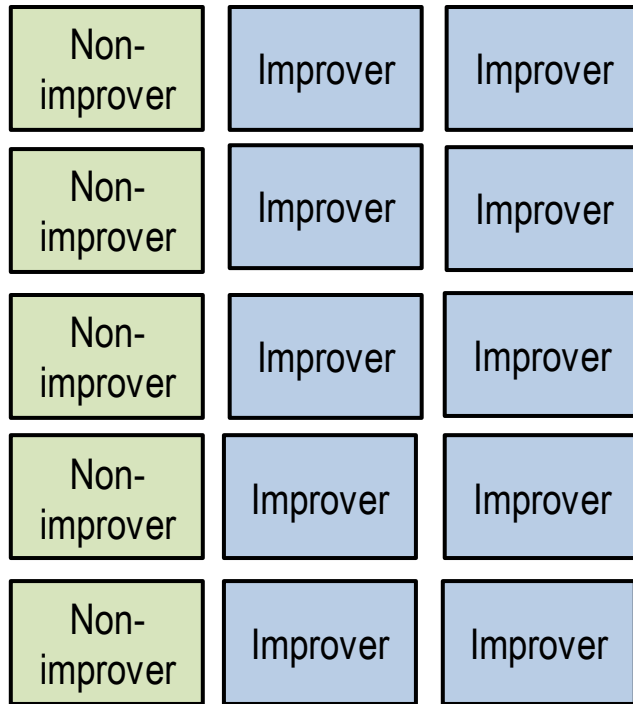
# Temperature and Heart Rate

- If there was no association between heart rate and body temperature, what is the probability we would get a correlation as high as 0.378 just by chance?
- If there is no association, we can break apart the temperatures and their corresponding heart rates. We will do this by shuffling one of the variables.

# Shuffling Cards

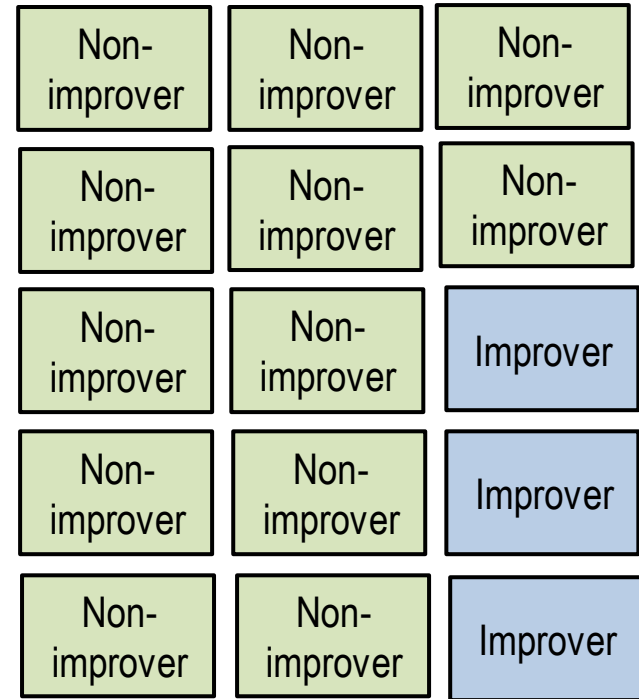
- Let's remind ourselves what we did with cards to find our simulated statistics.
- With two proportions, we wrote the response on the cards, shuffled the cards and placed them into two piles corresponding to the two categories of the explanatory variable.
- With two means we did the same thing except this time the responses were numbers instead of words.

# Dolphin Therapy



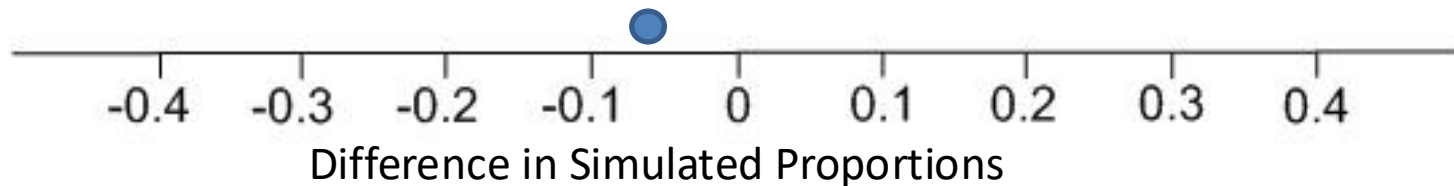
66.7% Improvers

# Control



33.3% Improvers

$$0.400 - 0.467 = -0.067$$



# Music

25.2	45.6
14.5	11.6
-7.0	18.6
12.6	12.1
34.5	30.5

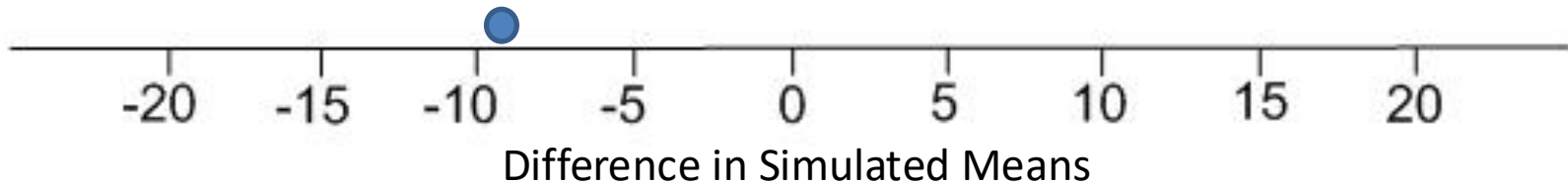
mean = 6.38

# No music

-10.7	-10.7	10.0
4.5	9.6	
2.2	2.4	
21.3	21.8	
-14.7	7.2	

mean = 16.12

$$6.38 - 16.12 = -9.74$$



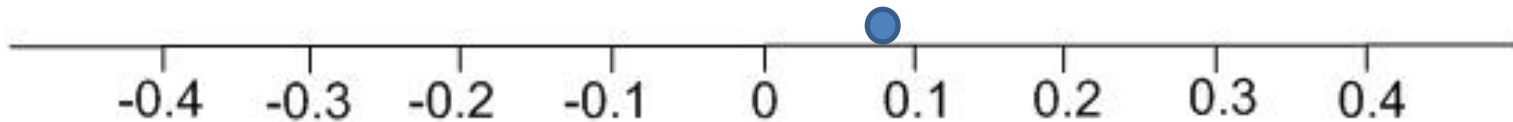
# Shuffling Cards

- Now how will this shuffling be different when both the response and the explanatory variable are quantitative?
- We can't put things in two piles anymore.
- We still shuffle values of the response variable, but this time place them next to two values of the explanatory variable.

# Body Temperature and Heart Rate

98.3°	98.2°	97.7°	98.5°	97.0°	98.8°	98.5°	98.7°	99.3°	97.8°
72	69	72	71	80	81	68	82	68	65
98.2°	99.9°	98.6°	98.6°	97.8°	98.4°	98.7°	97.4°	96.7°	98.0°
71	79	86	82	58	84	73	57	62	89

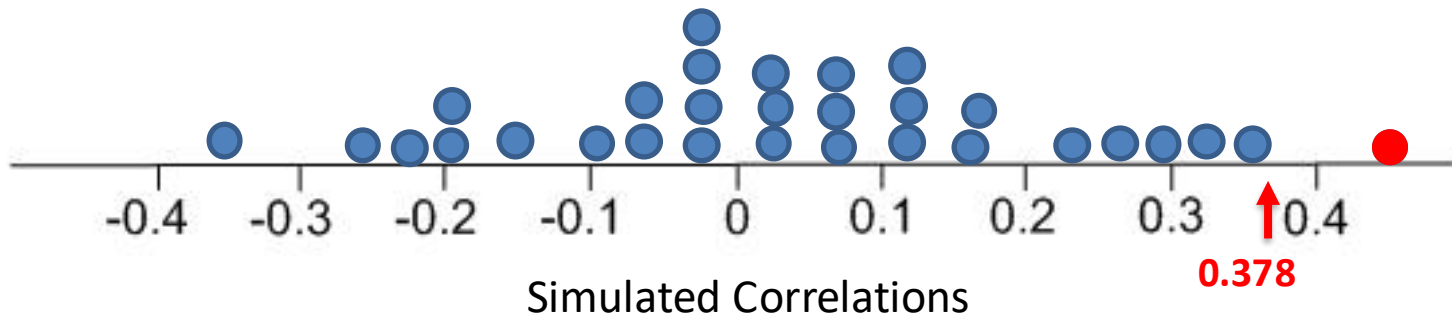
$r = 0.078$



Simulated Correlations

# More Simulations

Only one simulated statistic out of 30 was as large or larger than our observed correlation of 0.378, hence our p-value for this null distribution is  $1/30 \approx 0.03$ .

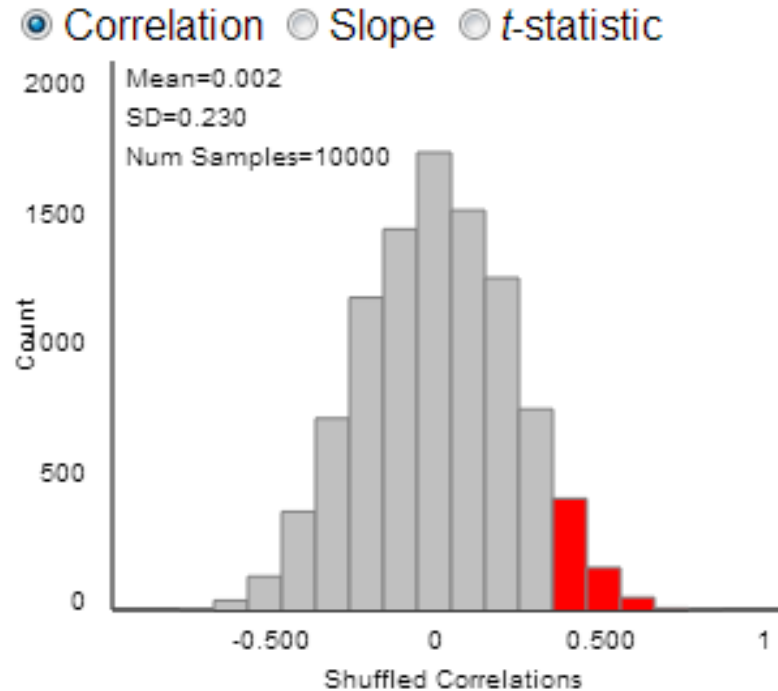


# Temperature and Heart Rate

- We can look at the output of 1000 shuffles with a distribution of 1000 simulated correlations.

# Temperature and Heart Rate

- Notice our null distribution is centered at 0 and somewhat symmetric.
- We found that 530/10000 times we had a simulated correlation greater than or equal to 0.378.



Count Samples

Count = 530/10000 (0.0530)

# Temperature and Heart Rate

- With a p-value of  $0.053 = 5.3\%$ , we almost but do not quite have statistical significance. We observe a positive linear association between body temperature and heart rate but this association is not statistically significant. Perhaps a larger sample should be investigated to get a better idea if the two variables are related or not.

### 3. Calculating correlation, r.

$\rho$  = rho = correlation of the population.

Suppose there are N people in the population,

X = temperature, Y = heart rate,

the mean and sd of temp in the pop. are  $\mu_x$  and  $\sigma_x$ ,

and the pop. mean and sd of heart rate are  $\mu_y$  and  $\sigma_y$ .

$$\rho = \frac{1}{N} \sum_{i=1}^N \left( \frac{x_i - \mu_x}{\sigma_x} \right) \left( \frac{y_i - \mu_y}{\sigma_y} \right).$$

Given a sample of size n, we estimate  $\rho$  using

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

This is in Appendix A.

# 4. Linear Regression

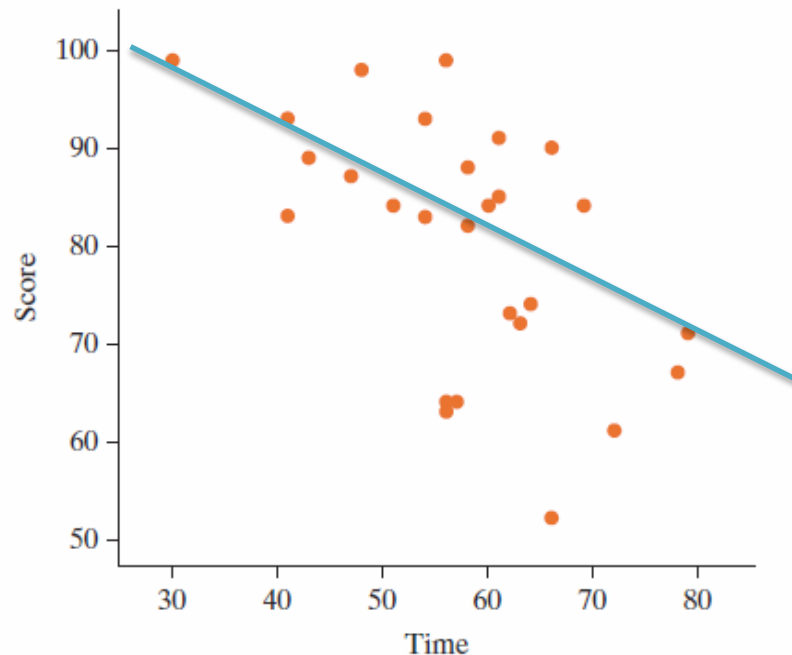
Section 10.3

# Introduction

- If we decide an association is linear, it is helpful to develop a mathematical model of that association.
- Helps make predictions about the response variable.
- The *least-squares regression line* is the most common way of doing this.

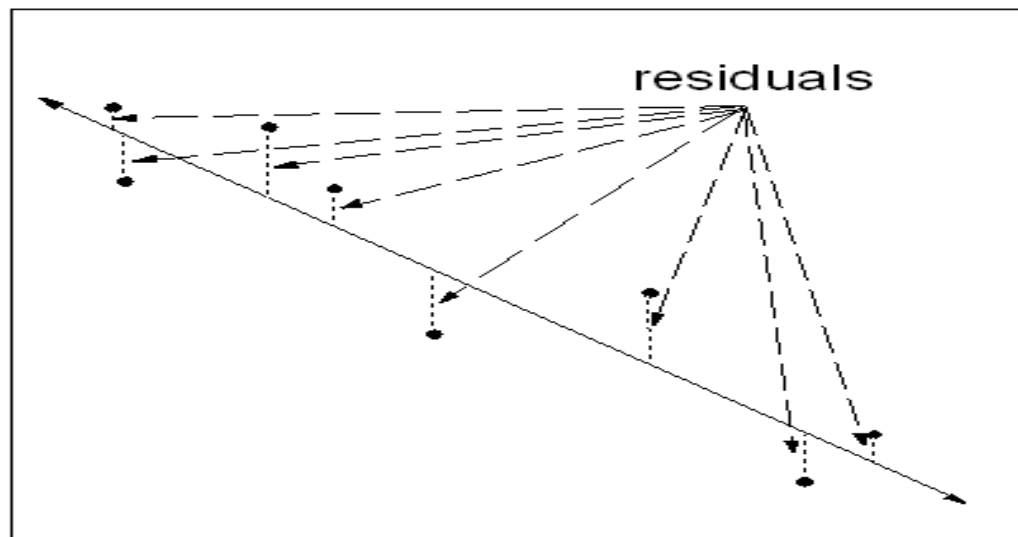
# Introduction

- Unless the points are perfectly linearly aligned, there will not be a single line that goes through every point.



# Introduction

- We want a line that minimizes the vertical distances between the line and the points
  - These distances are called **residuals**.
  - The line we will find actually minimizes the sum of the squares of the residuals.
  - This is called a **least-squares regression line**.

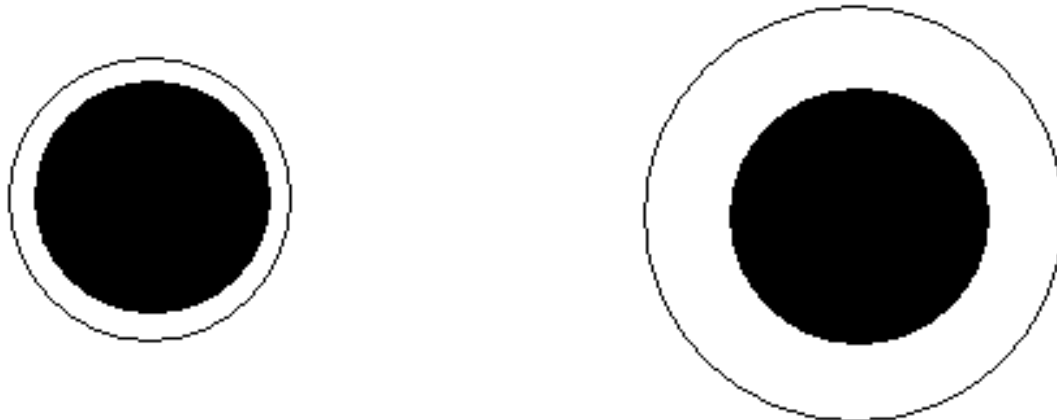


# Are Dinner Plates Getting Larger?

*Example 10.3*

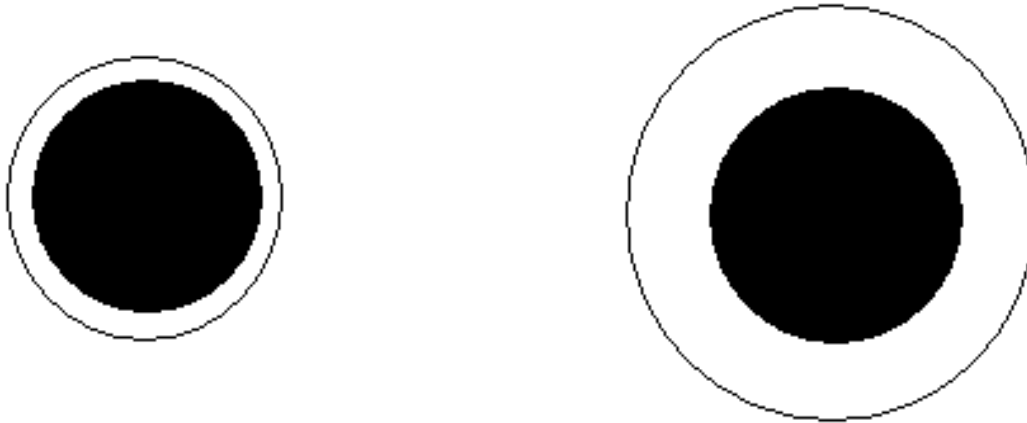
# Growing Plates?

- There are many recent articles and TV reports about the obesity problem.
- One reason some have given is that the size of dinner plates are increasing.
- Are these black circles the same size, or is one larger than the other?



# Growing Plates?

- They appear to be the same size for many, but the one on the right is about 20% larger than the left.



- This suggests that people will put more food on larger dinner plates without knowing it.
- There is name for this phenomenon: *Delboeuf illusion*.

# Growing Plates?

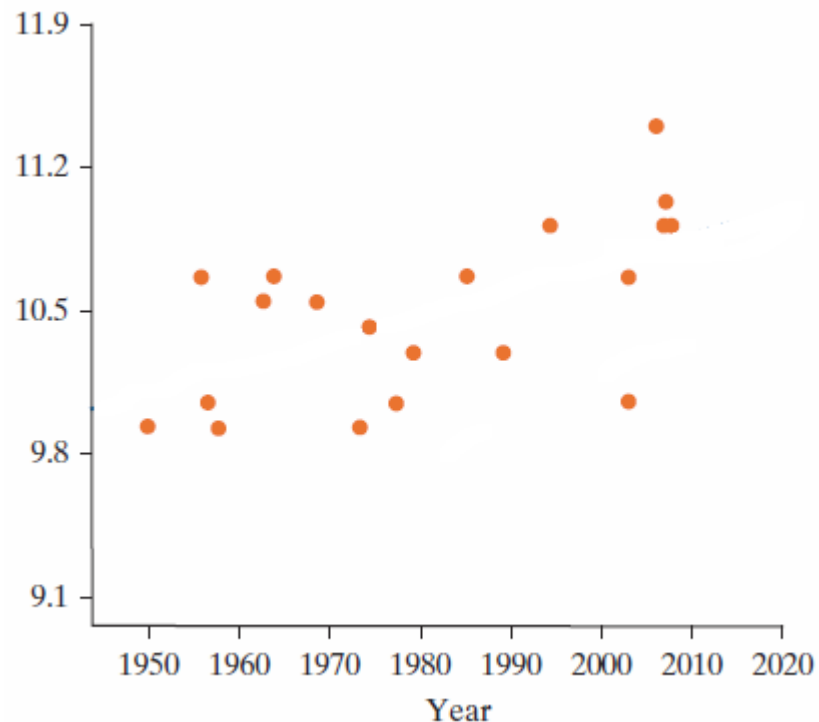
- Researchers gathered data to investigate the claim that dinner plates are growing
- American dinner plates sold on ebay on March 30, 2010 (Van Ittersum and Wansink, 2011)
- Year manufactured and diameter are given.

**TABLE 10.1** Data for size (diameter, in inches) and year of manufacture for 20 American-made dinner plates

<b>Year</b>	1950	1956	1957	1958	1963	1964	1969	1974	1975	1978
<b>Size</b>	10	10.75	10.125	10	10.625	10.75	10.625	10	10.5	10.125
<b>Year</b>	1980	1986	1990	1995	2004	2004	2007	2008	2008	2009
<b>Size</b>	10.375	10.75	10.375	11	10.75	10.125	11.5	11	11.125	11

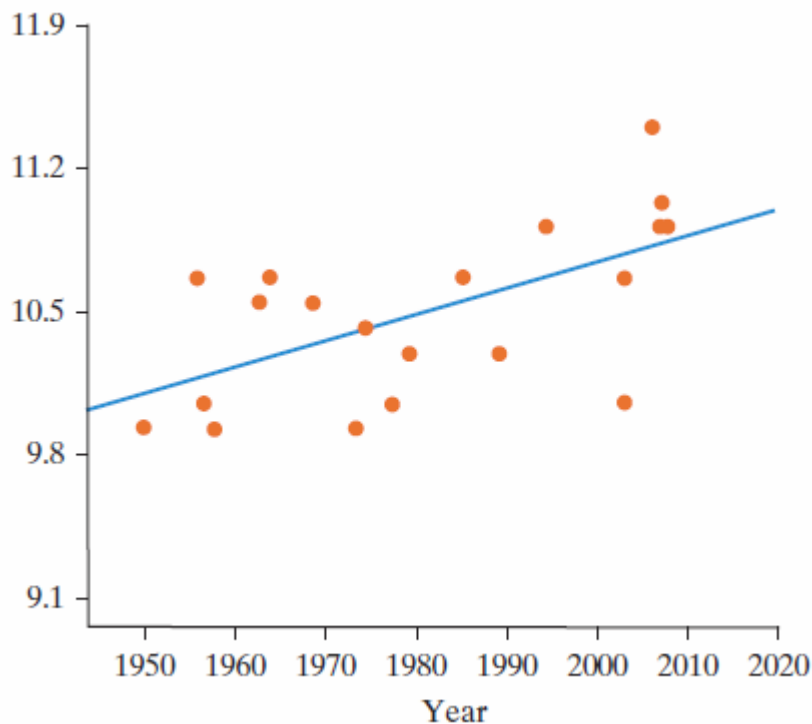
# Growing Plates?

- Both year (explanatory variable) and diameter in inches (response variable) are quantitative.
- Each dot in this scatterplot represents one plate.



# Growing Plates?

- The association appears to be roughly linear.
- The least squares regression line is added.
- The line slopes upward, but is the slope significant?



# Regression Line

The regression equation is  $\hat{y} = a + bx$ :

- $a$  is the  $y$ -intercept
  - $b$  is the slope
  - $x$  is a value of the explanatory variable
  - $\hat{y}$  is the predicted value for the response variable
- For a specific value of  $x$ , the corresponding distance  $y - \hat{y}$  (or actual – predicted) is a residual

# Regression Line

- The least squares line for the dinner plate data is  $\hat{y} = -14.8 + 0.0128x$
- Or  $\widehat{\text{diameter}} = -14.8 + 0.0128(\text{year})$
- This allows us to predict plate diameter for a particular year.

# Slope

$$\hat{y} = -14.8 + 0.0128x$$

- What is the predicted diameter for a plate manufactured in 2000?
  - $-14.8 + 0.0128(2000) = 10.8$  in.
- What is the predicted diameter for a plate manufactured in 2001?
  - $-14.8 + 0.0128(2001) = 10.8128$  in.
- How does this compare to our prediction for the year 2000?
  - 0.0128 larger
- Slope  $b = 0.0128$  means that diameters are predicted to increase by 0.0128 inches per year on average

# Slope

- Slope is the predicted change in the response variable for one-unit change in the explanatory variable.
- Both the slope and the correlation coefficient for this study were positive.
  - The slope is 0.0128
  - The correlation is 0.604
- The slope and correlation coefficient will always have the same sign.

# Slope of regression line.

- Suppose  $\hat{y} = a + bx$  is the regression line.
- The slope  $b$  of the regression line is  $b = r \frac{s_y}{s_x}$ .

This is usually the thing of primary interest to interpret, as the predicted increase in  $y$  for every unit increase in  $x$ .

- Beware of assuming causation though, esp. with observational studies. Be wary of extrapolation too.
- The intercept  $a = \bar{y} - b \bar{x}$ .
- The SD of the residuals is  $\sqrt{1 - r^2} s_y$ .  
This is a good estimate of how much the regression predictions will typically be off by.

# y-intercept

- The y-intercept is where the regression line crosses the y-axis. It is the predicted response when the explanatory variable equals 0.
- We had a y-intercept of -14.8 in the dinner plate equation. What does this tell us about our dinner plate example?
  - Dinner plates in year 0 would be predicted to be -14.8 inches???
- How can it be negative?
  - The equation works well within the range of values given for the explanatory variable, but fails outside that range.
- Our equation should only be used to predict the size of dinner plates from about 1950 to 2010.