

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

1. Calculating correlation.
2. Testing correlation.
3. Linear regression.
4. Slope of regression line.
5. Goodness of fit.

No class Thu Nov 24, Thanksgiving.

Read ch10.

Hw4 is 10.1.8, 10.3.14, 10.3.21, 10.4.11 and is due Tue Nov 29.

The final Fri Dec 9, 8am-11, right here, will be on ch1-10.

Bring a PENCIL and CALCULATOR and any books or notes you want. No computers.

<http://www.stat.ucla.edu/~frederic/13/F16> .

## 1. 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.

# 2. 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. ( $\rho > 0$ )

# 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

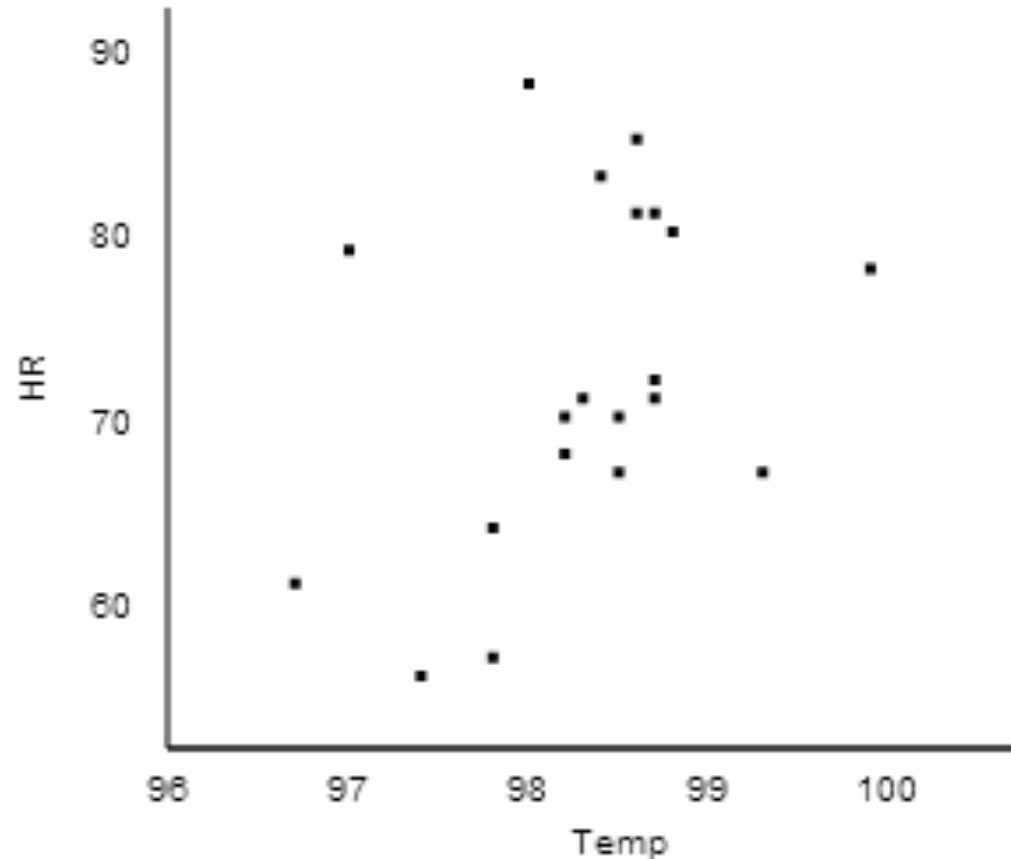
Collect the Data

Tmp	98.3	98.2	98.7	98.5	97.0	98.8	98.5	98.7	99.3	97.8
HR	72	69	72	71	80	81	68	82	68	65
Tmp	98.2	99.9	98.6	98.6	97.8	98.4	98.7	97.4	96.7	98.0
HR	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

Non-improver	Improver	Improver
Non-improver	Improver	Improver
Non-improver	Improver	Improver
Non-improver	Improver	Improver
Non-improver	Improver	Improver

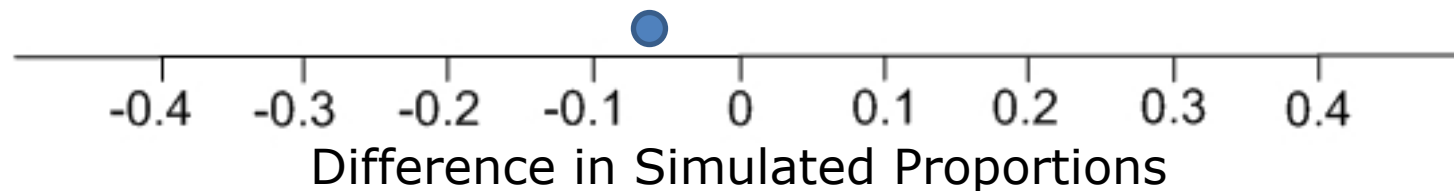
40.0%  
Improvers

# Control

Non-improver	Non-improver	Non-improver
Non-improver	Non-improver	Non-improver
Non-improver	Non-improver	Improver
Non-improver	Non-improver	Improver
Non-improver	Non-improver	Improver

20.0%  
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

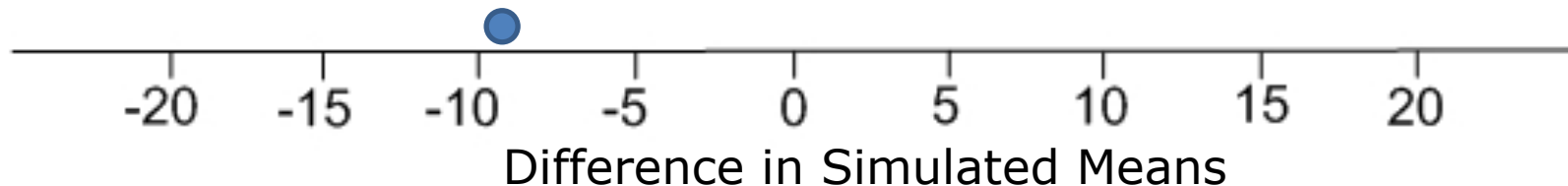
mean = 6.38

## 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

$$6.38 - 16.12 = -9.74$$



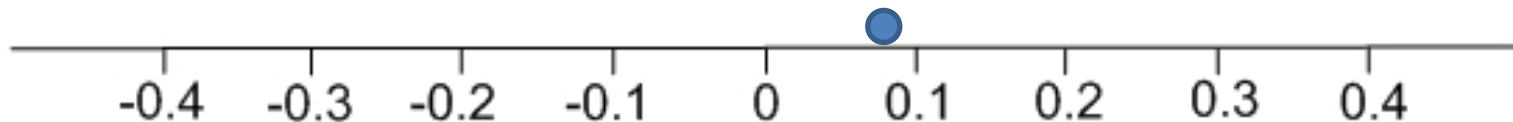
# Shuffling Cards

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

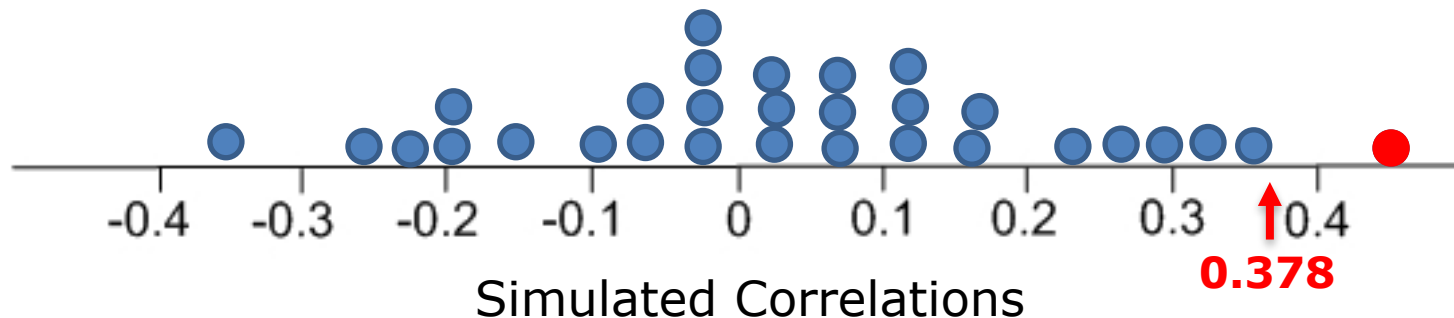
$r = 0.078$



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  $1/30 \approx 0.03$ .



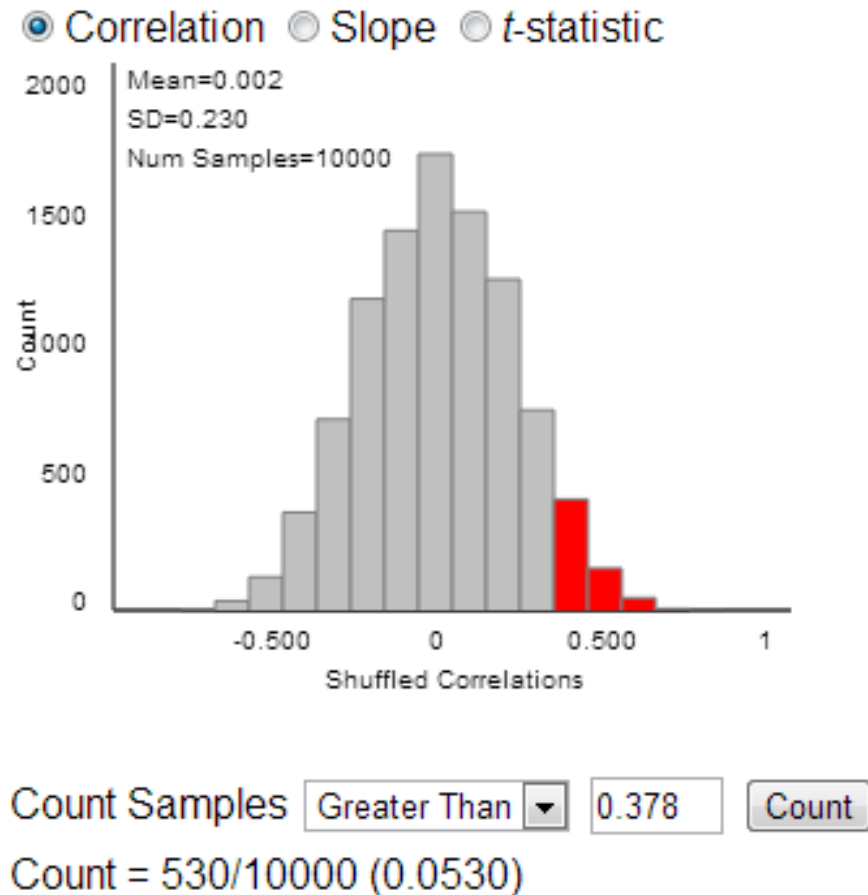
# Temperature and Heart Rate

- 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.



# Temperature and Heart Rate

- With a p-value of  $0.053 = 5.3\%$ , we almost but do not quite have statistical significance. This is moderate evidence of a positive linear association between body temperature and heart rate. Perhaps a larger sample would give a smaller p-value.

# 3. Least Squares Regression

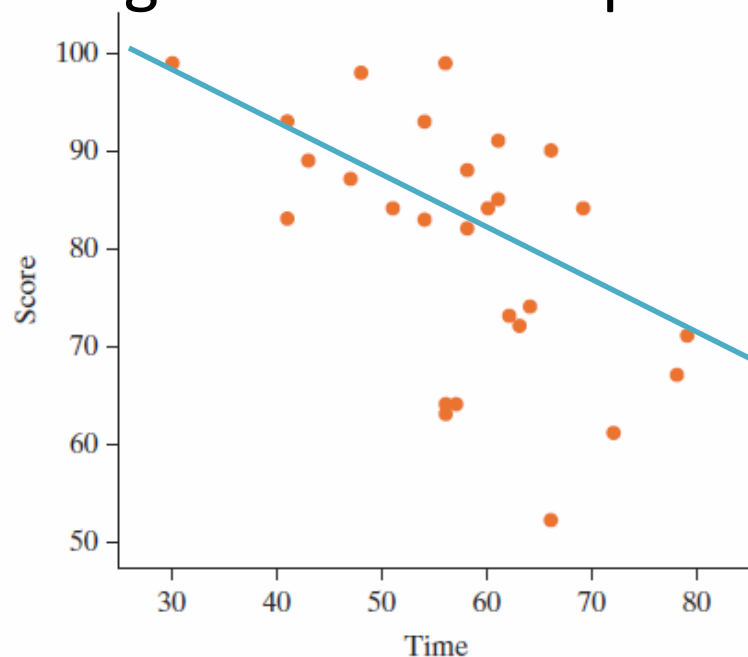
Section 10.3

# 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.

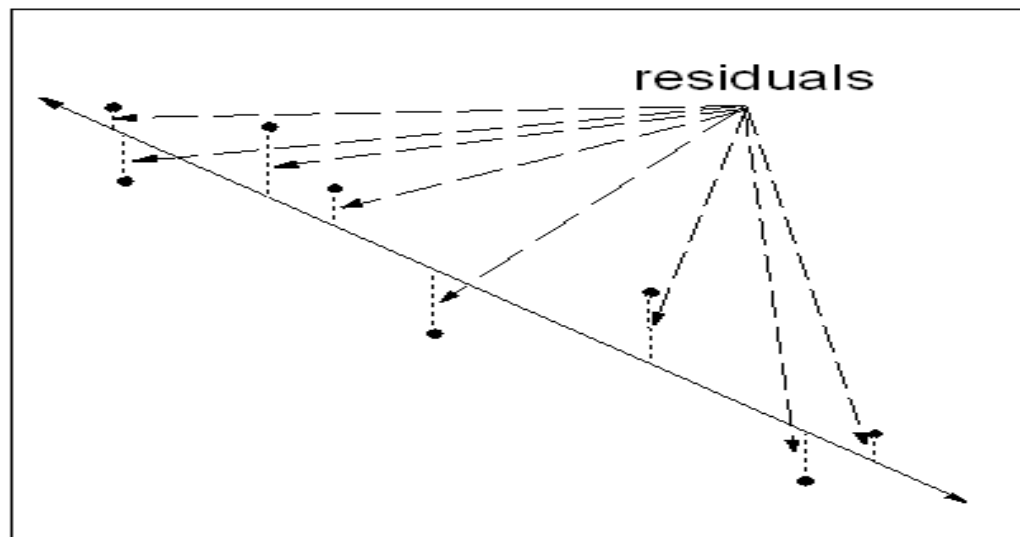
# Introduction

- Unless the points are perfectly linearly aligned, there will not be a single line that goes through every point.
- We want a line that gets as close as possible to all the points.



# 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**.

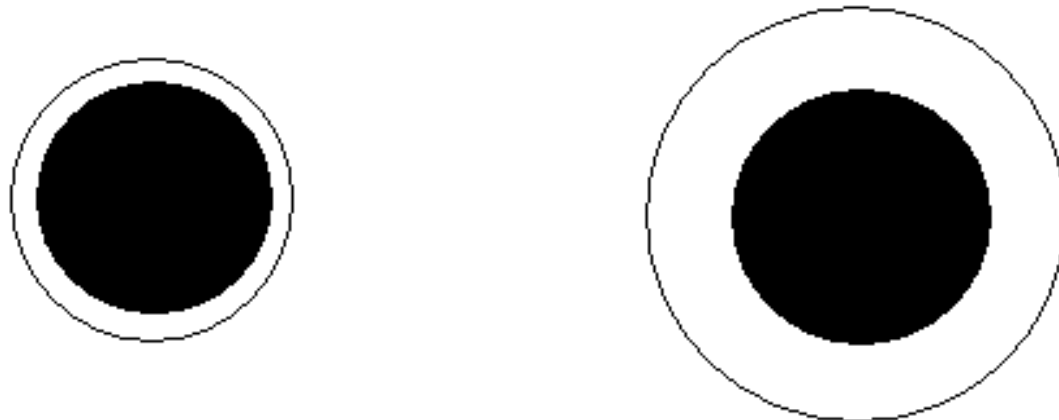


# Are Dinner Plates Getting Larger?

*Example 10.3*

# Growing Plates?

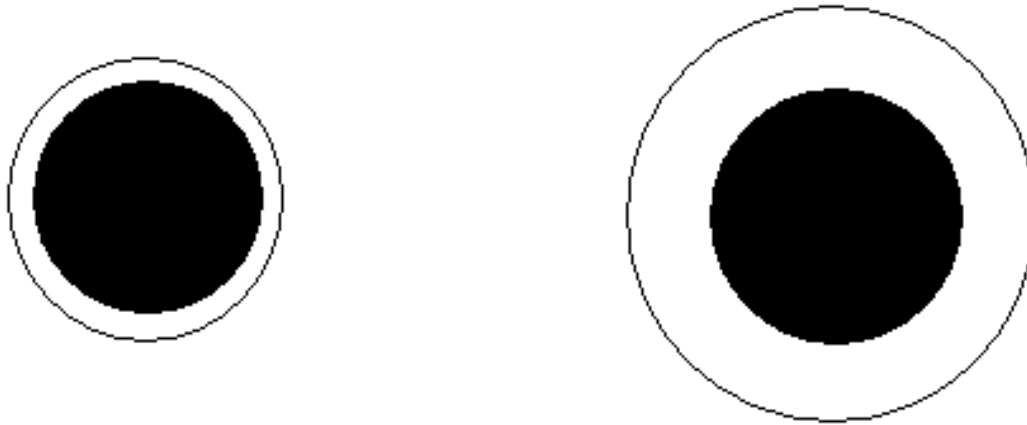
- There are many recent articles and TV reports about the obesity problem.
- One reason some have given is that the size of dinner plates are increasing.
- Are these black circles the same size, or is one larger than the other?





# 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?

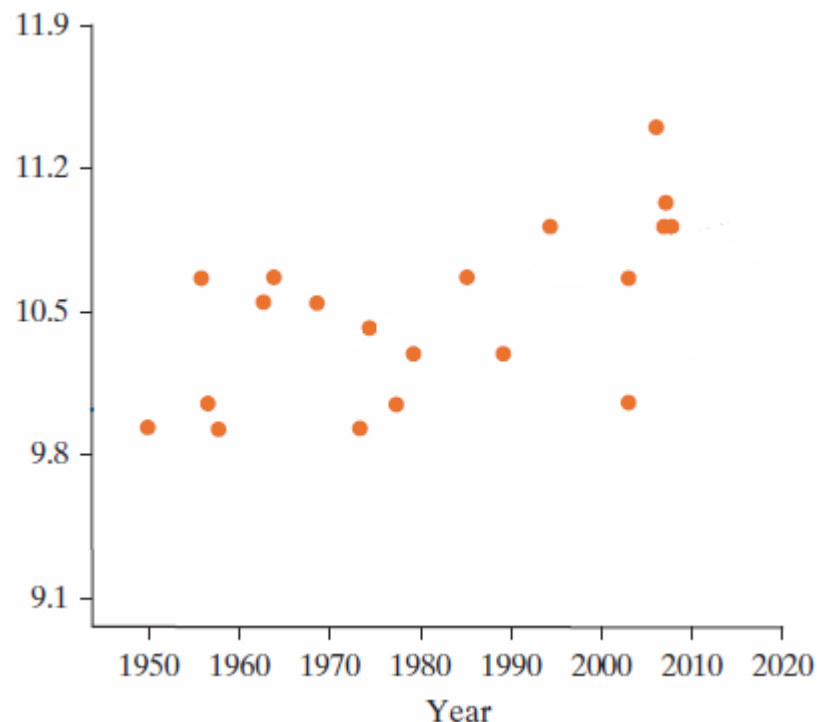
- Researchers gathered data to investigate the claim that dinner plates are growing
- American dinner plates sold on ebay on March 30, 2010 (Van Ittersum and Wansink, 2011)
- Year manufactured and diameter are given.

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

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

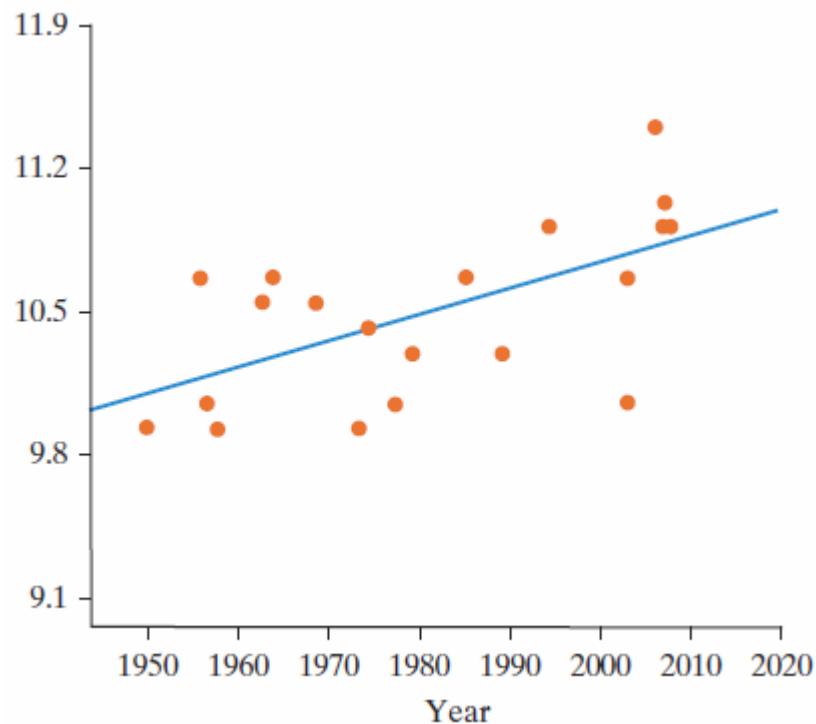
# Growing Plates?

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



# Growing Plates?

- The association appears to be roughly linear
- The least squares regression line is added
- How can we describe this line?



# 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.

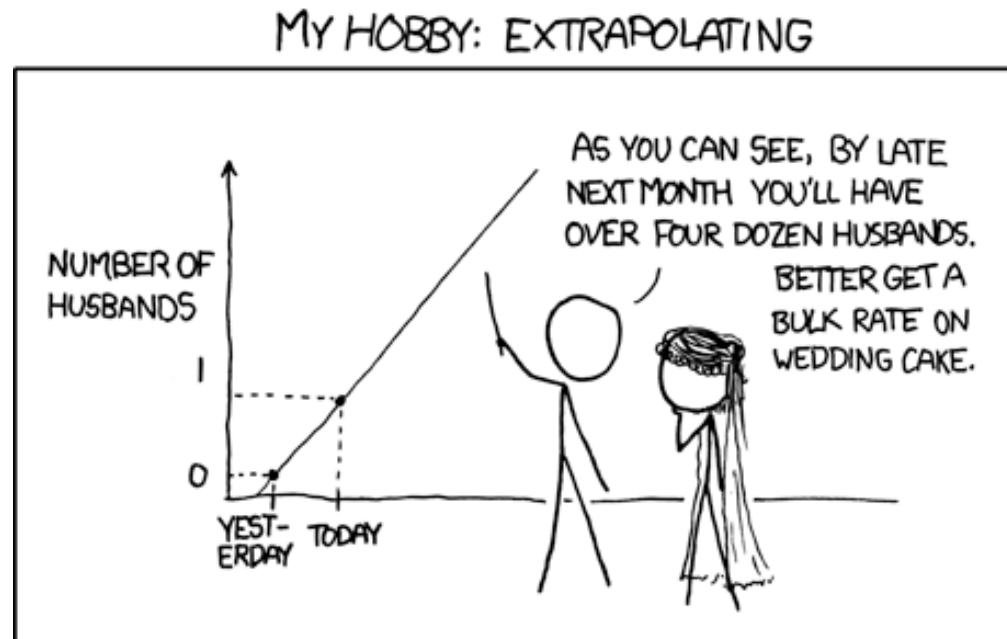


# $y$ -intercept

- The  $y$ -intercept is where the regression line crosses the  $y$ -axis or 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 were -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.

# Extrapolation

- Predicting values for the response variable for values of the explanatory variable that are outside of the range of the original data is called ***extrapolation***.



# Coefficient of Determination

- While the intercept and slope have meaning in the context of year and diameter, remember that the correlation does not. It is just 0.604.
- However, the square of the correlation (coefficient of determination or  $r^2$ ) does have meaning.
- $r^2 = 0.604^2 = 0.365$  or 36.5%
- 36.5% of the variation in plate size (the response variable) can be explained by its linear association with the year (the explanatory variable).

# Learning Objectives for Section 10.3

- Understand that one way a scatterplot can be summarized is by fitting the best-fit (least squares regression) line.
- Be able to interpret both the slope and intercept of a best-fit line in the context of the two variables on the scatterplot.
- Find the predicted value of the response variable for a given value of the explanatory variable.
- Understand the concept of residual and find and interpret the residual for an observational unit given the raw data and the equation of the best fit (regression) line.
- Understand the relationship between residuals and strength of association and that the best-fit (regression) line this minimizes the sum of the squared residuals.

# Learning Objectives for Section 10.3

- Find and interpret the coefficient of determination ( $r^2$ ) as the squared correlation and as the percent of total variation in the response variable that is accounted for by the linear association with the explanatory variable.
- Understand that extrapolation is when a regression line is used to predict values outside of the range of observed values for the explanatory variable.
- Understand that when slope = 0 means no association, slope < 0 means negative association, slope > 0 means positive association, and that the sign of the slope will be the same as the sign of the correlation coefficient.
- Understand that influential points can substantially change the equation of the best-fit line.

## 4. Slope of regression line.

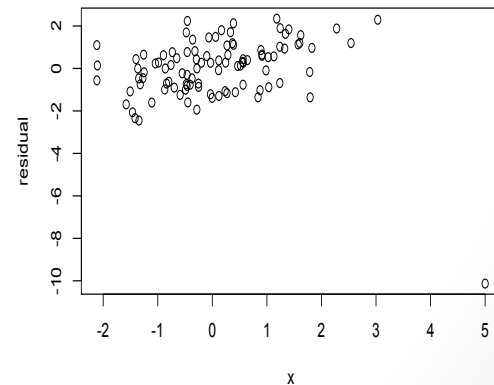
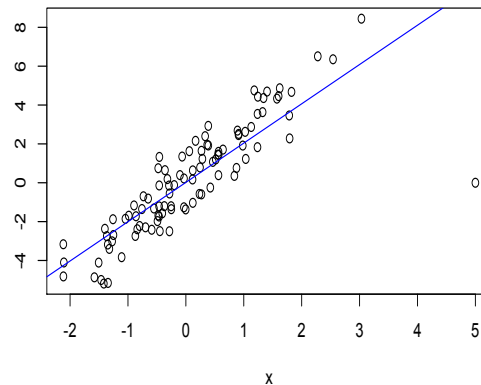
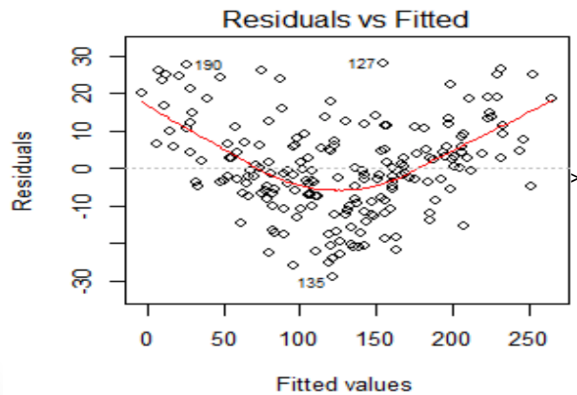
- Suppose  $\hat{y} = a + bx$  is the regression line.
- The slope  $b$  of the regression line is  $b = r \frac{s_y}{s_x}$ .

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

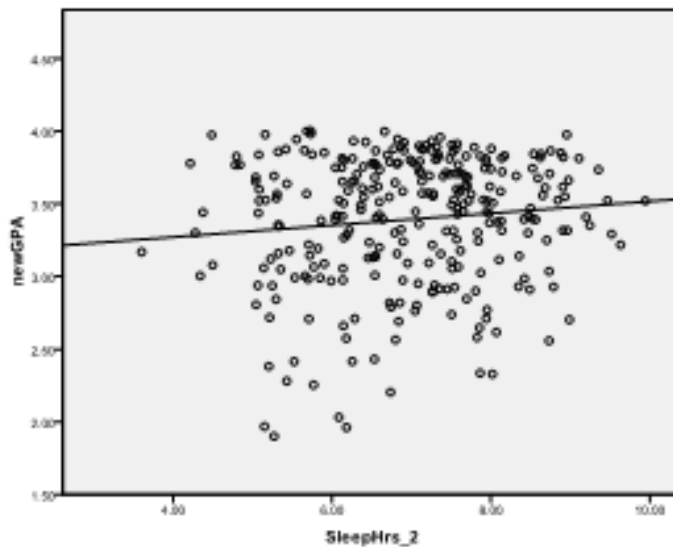
- Beware of assuming causation though, esp. with observational studies. Be wary of extrapolation too.
- The intercept  $a = \bar{y} - b \bar{x}$ .
- The SD of the residuals is  $\sqrt{1 - r^2} s_y$ .  
This is a good estimate of how much the regression predictions will typically be off by.

# 5. How well does the line fit?

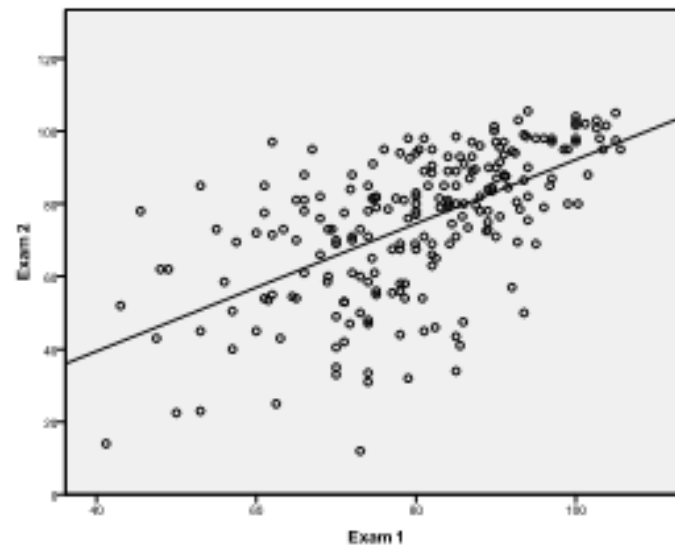
- $r^2$  is a measure of fit. It indicates the amount of scatter around the best fitting line.
- Residual plots can indicate curvature, outliers, or heteroskedasticity.
- $\sqrt{1 - r^2} s_y$  is useful as a measure of how far off predictions would have been on average.



- Heteroskedasticity: when the variability in  $y$  is not constant as  $x$  varies.



(a)



(b)



## 5. How well does the line fit?

- $r^2$  is a measure of fit. It indicates the amount of scatter around the best fitting line.
- Residual plots can indicate curvature, outliers, or heteroskedasticity.
- $\sqrt{1 - r^2} s_y$  is useful as a measure of how far off predictions would have been on average.

Is the estimated slope  $b$  significantly different from 0? Is the correlation  $r$  significantly different from 0? These are really the same test. We will discuss testing this next time.