

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

1. Regression line, continued.
2. Calculating correlation.
3. Slope of regression line.
4. Goodness of fit.
5. Common problems with regression.
6. Testing the slope or correlation.

No class Thu Nov 23, Thanksgiving.

Read ch10.

Hw4 is due Tue Nov 28.

The final Mon, Dec 11, 3-6pm.

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

<http://www.stat.ucla.edu/~frederic/13/F17>.

# 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

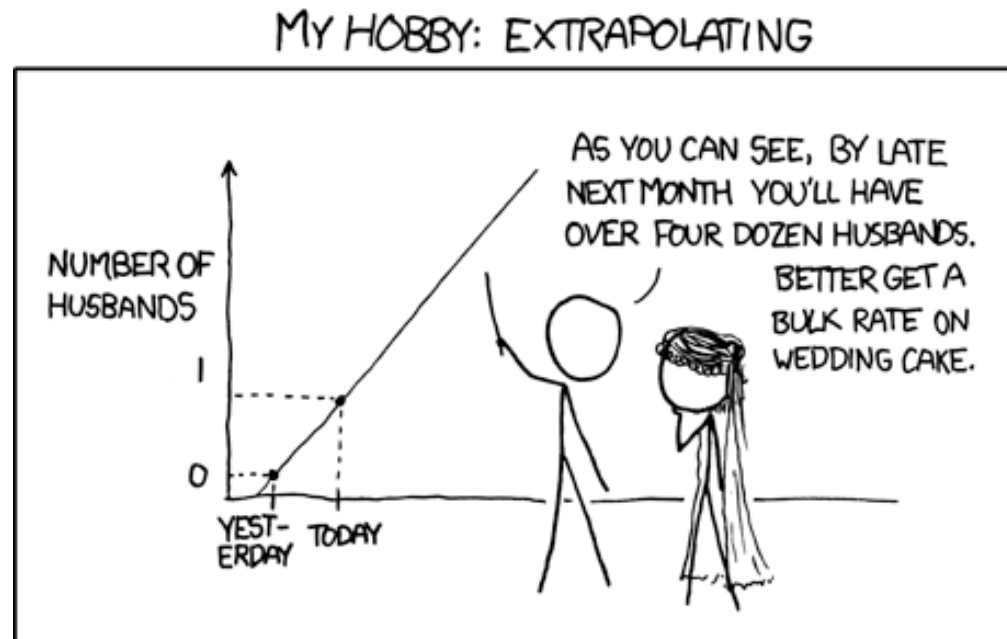
- Slope is the predicted change in the response variable for one-unit change in the explanatory variable.
- Both the slope and the correlation coefficient for this study were positive.
  - The slope is 0.0128
  - The correlation is 0.604
- The slope and correlation coefficient will always have the same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

## 2. 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.



### 3. Slope of regression line.

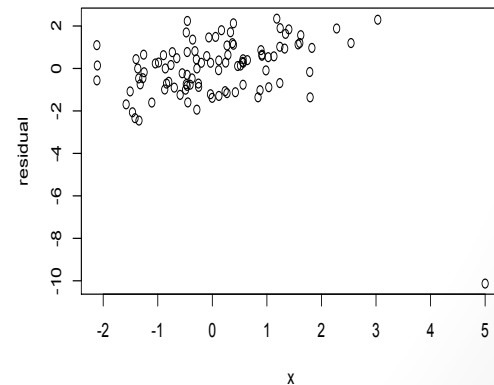
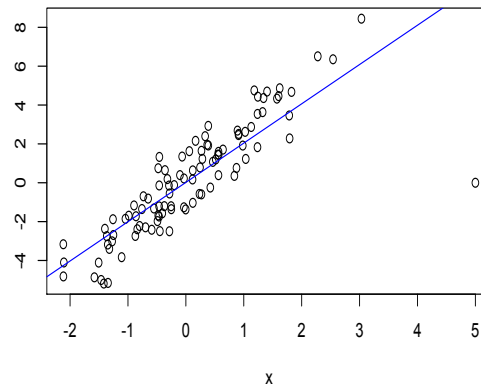
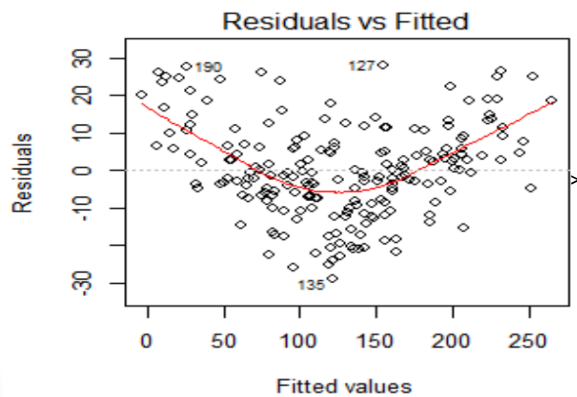
- Suppose  $\hat{y} = a + bx$  is the regression line.
- The slope  $b$  of the regression line is  $b = r \frac{s_y}{s_x}$ .

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

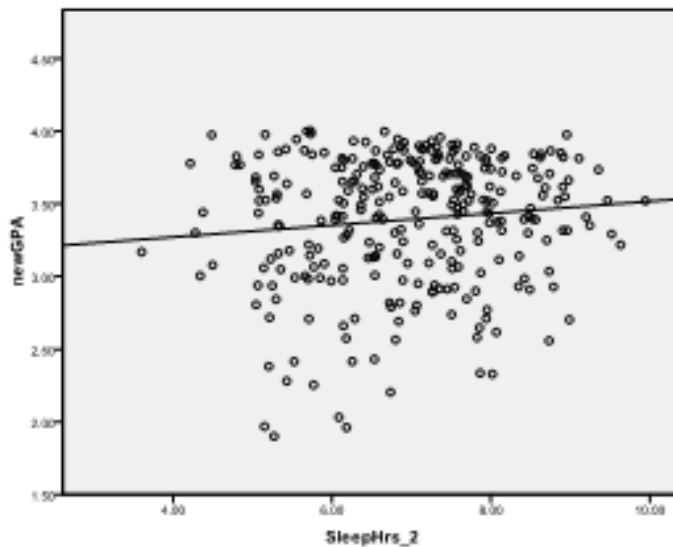
- Beware of assuming causation though, esp. with observational studies. Be wary of extrapolation too.
- The intercept  $a = \bar{y} - b \bar{x}$ .
- The SD of the residuals is  $\sqrt{1 - r^2} s_y$ .  
This is a good estimate of how much the regression predictions will typically be off by.

## 4. How well does the line fit?

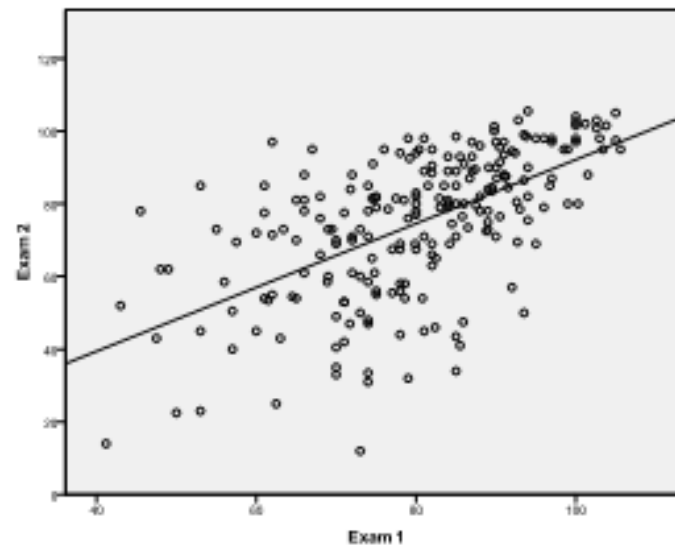
- $r^2$  is a measure of fit. It indicates the amount of scatter around the best fitting line.
- Residual plots can indicate curvature, outliers, or heteroskedasticity.
- $\sqrt{1 - r^2} s_y$  is useful as a measure of how far off predictions would have been on average.



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



(a)



(b)

# 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.

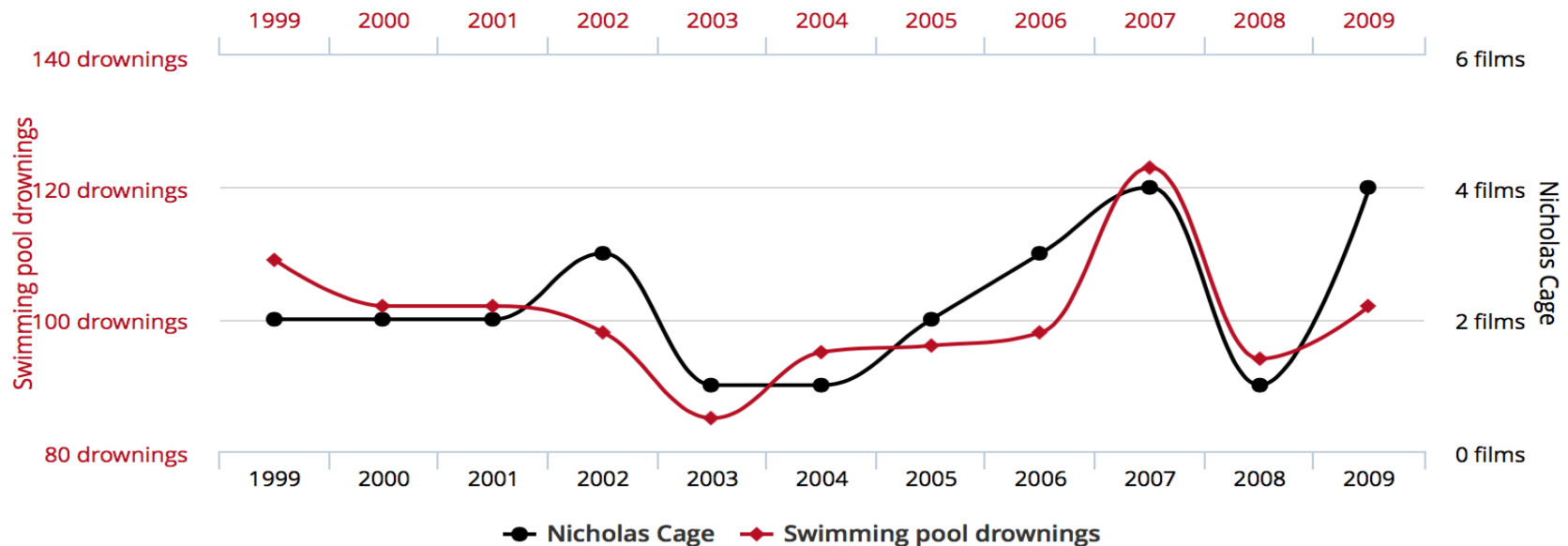
Note that regression residuals have mean zero, whether the line fits well or poorly.

# 5. Common problems with regression.

- a. Correlation is not causation.  
Especially with observational data.

**Number of people who drowned by falling into a pool**  
correlates with  
**Films Nicolas Cage appeared in**

Correlation: 66.6% ( $r=0.666004$ )



tylervigen.com

# Common problems with regression.



# Common problems with regression.

Holmes and Willett (2004) reviewed all prospective studies on fat consumption and breast cancer with at least 200 cases of breast cancer. "Not one study reported a significant positive association with total fat intake.... Overall, no association was observed between intake of total, saturated, monounsaturated, or polyunsaturated fat and risk for breast cancer."

They also state "The dietary fat hypothesis is largely based on the observation that national per capita fat consumption is highly correlated with breast cancer mortality rates. However, per capita fat consumption is highly correlated with economic development. Also, low parity and late age at first birth, greater body fat, and lower levels of physical activity are more prevalent in Western countries, and would be expected to confound the association with dietary fat."

# Common problems with regression.

- b. Extrapolation.

If the birthrate remains at **1.19** children per woman, South Korea could face natural extinction by **2750.**

Source:  
<http://blogs.wsj.com/korearealtime/2014/08/26/south-korea-birthrate-hits-lowest-on-record/>

BROOKINGS



# Common problems with regression.

- b. Extrapolation.
- Often researchers extrapolate from high doses to low.

D.M. Odom et al.

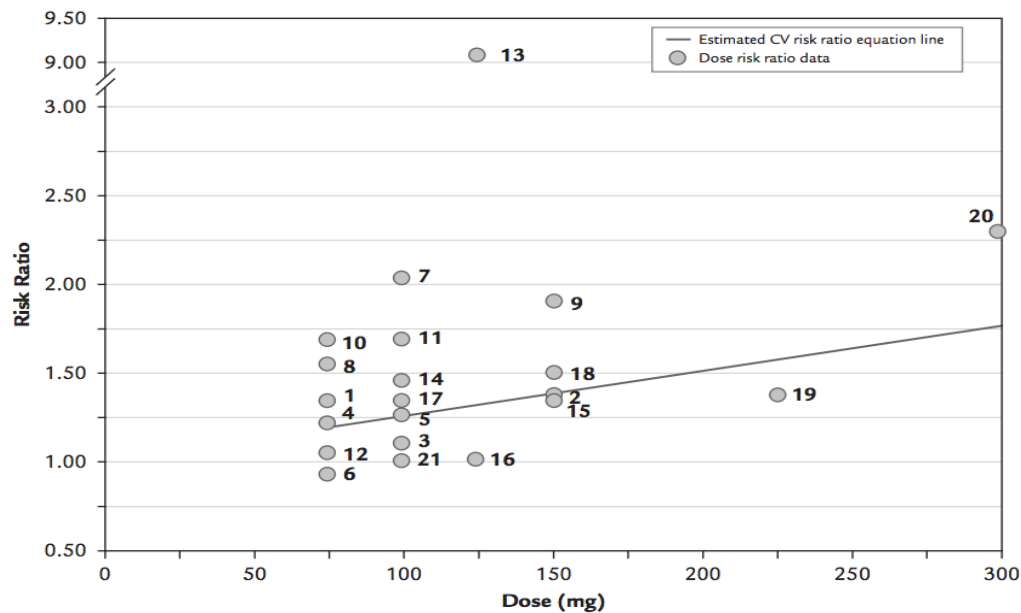


Figure 4. Relationship between diclofenac daily dose and the estimated risk ratio of a cardiovascular event. Numbers correspond to the observations in [Table III](#).

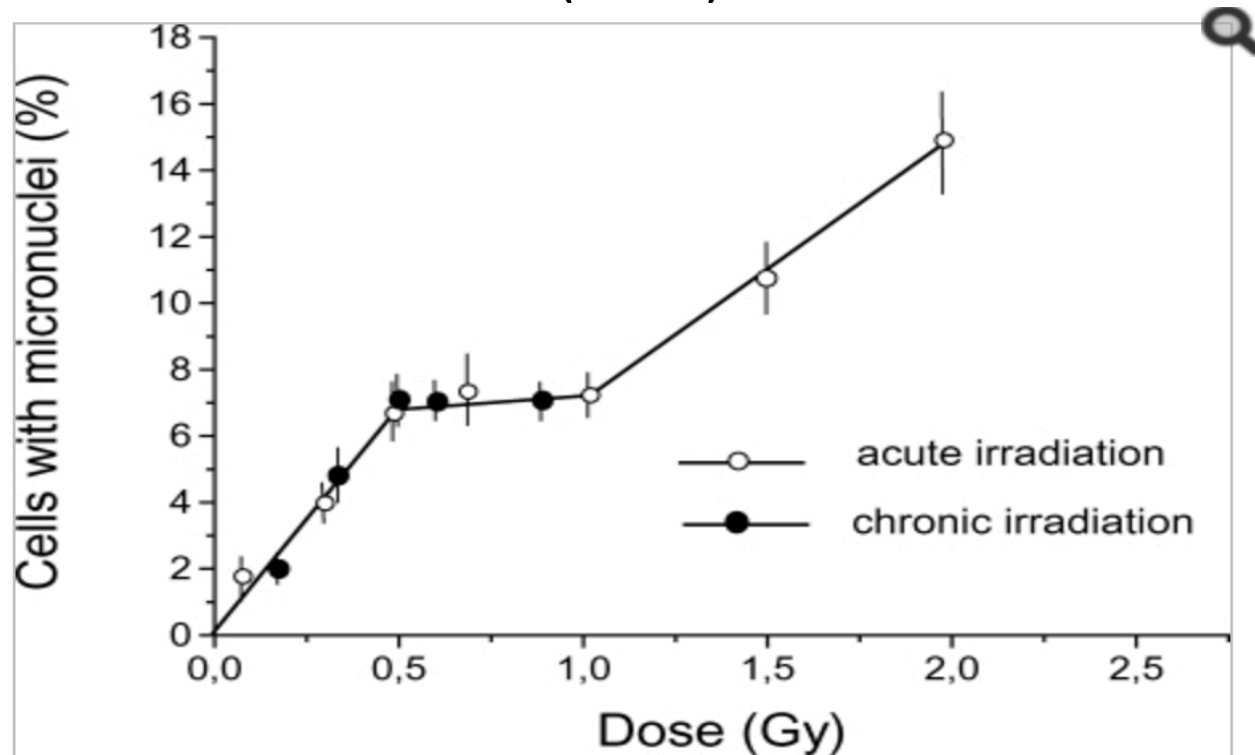
# Common problems with regression.

- b. Extrapolation.

The relationship can be nonlinear though.

Researchers also often extrapolate from animals to humans.

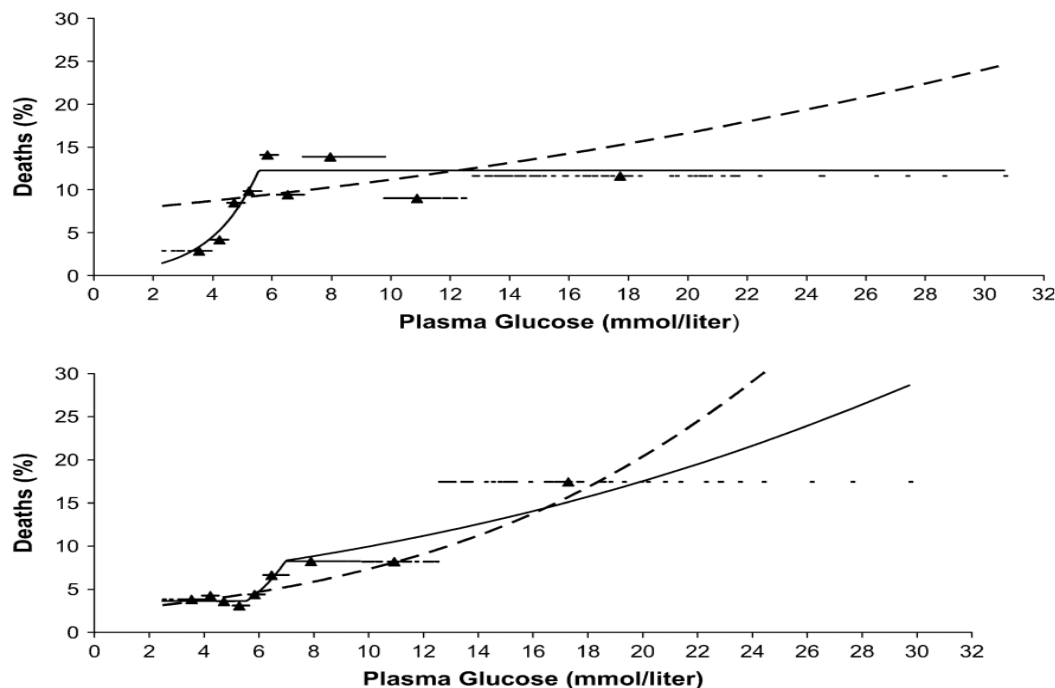
Zaichkina et al. (2004) on hamsters



# Common problems with regression.

- c. Curvature.

The best fitting line might fit poorly. Port et al. (2005).

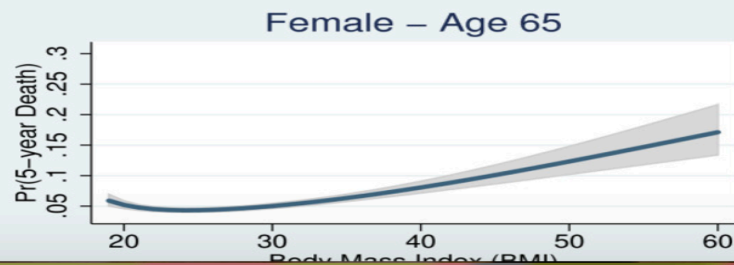
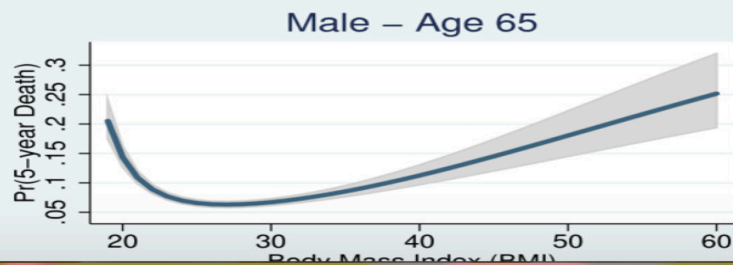
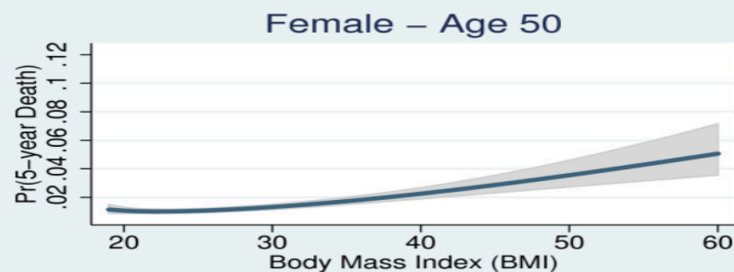
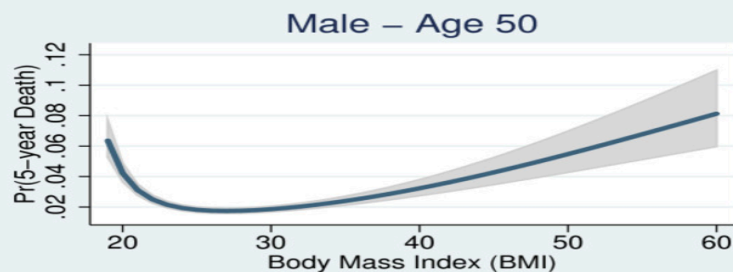
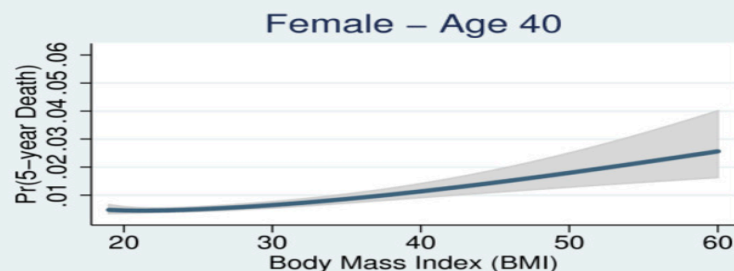
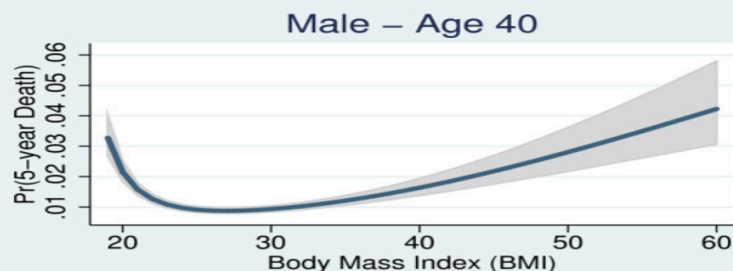


**FIGURE 4.** Adjusted 2-year rates of death from all causes for men (upper panel) and women (lower panel) separately, by glucose level, predicted by three models, Framingham Heart Study, 1948–1978. Linear model (dashed curve); optimal spline models (solid curve). The horizontal dashed

# Common problems with regression.

- c. Curvature.

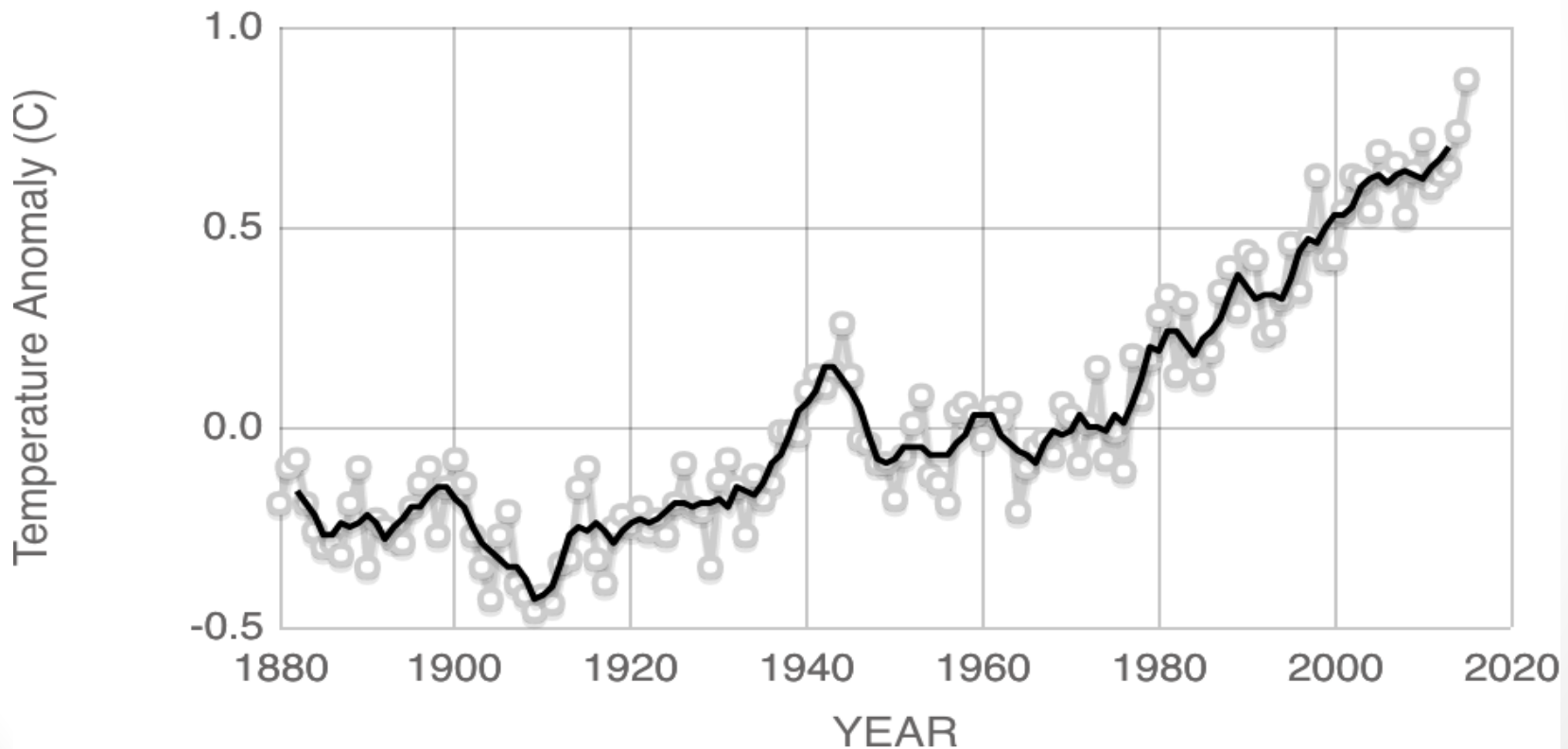
The best fitting line might fit poorly. Wong et al. (2011).



# Common problems with regression.

- d. Statistical significance.

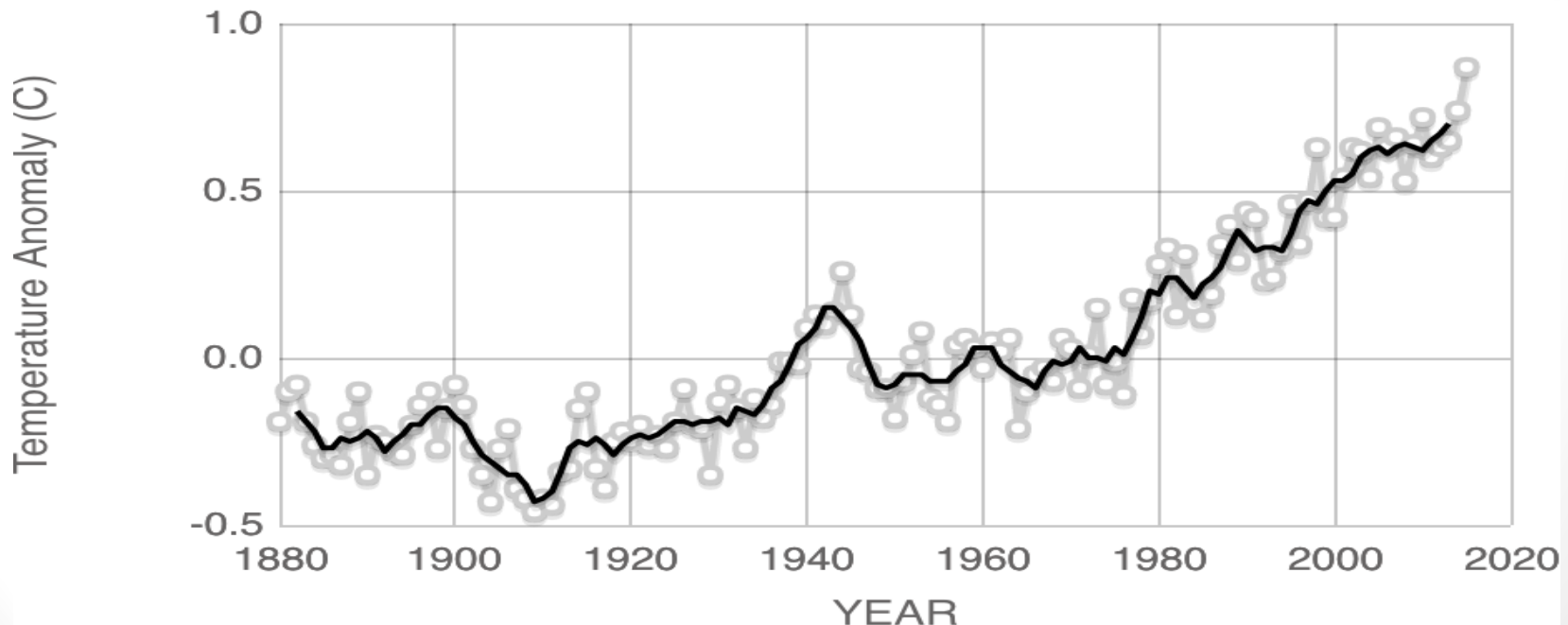
Could the observed correlation just be due to chance alone?



# Common problems with regression.

- d. Statistical significance.

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 it out.



# 6. Inference for the Regression Slope: Theory-Based Approach

Section 10.5

Do students who spend more time  
in non-academic activities tend to  
have lower GPAs?

Example 10.4



# Do students who spend more time in non-academic activities tend to have lower GPAs?

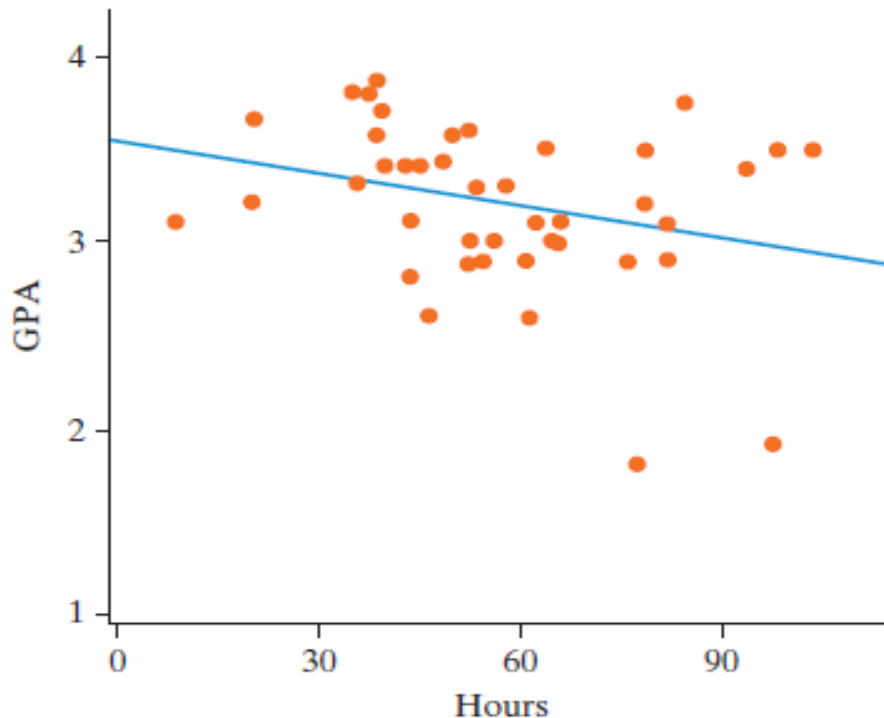
- The subjects were 34 undergraduate students from the University of Minnesota.
- They were asked questions about how much time they spent in activities like work, watching TV, exercising, non-academic computer use, etc. as well as what their current GPA was.
- We are going to test to see if there is a **negative** association between the number of hours per week spent on nonacademic activities and GPA.

# Hypotheses

- Null Hypothesis: There is no association between the number of hours students spend on nonacademic activities and student GPA in the population.
- Alternative Hypothesis: There is a negative association between the number of hours students spend on nonacademic activities and student GPA in the population.

# Descriptive Statistics

- $\widehat{GPA} = 3.60 - 0.0059(\text{nonacademic hours})$ .
- What do the slope and y-intercept mean?



# Shuffle to Develop Null Distribution

- We are going to shuffle just as we did with correlation to develop a null distribution.
- The only difference is that we will be calculating the slope each time and using that as our statistic.
- a test of association based on slope is equivalent to a test of association based on a correlation coefficient.

# Beta vs Rho

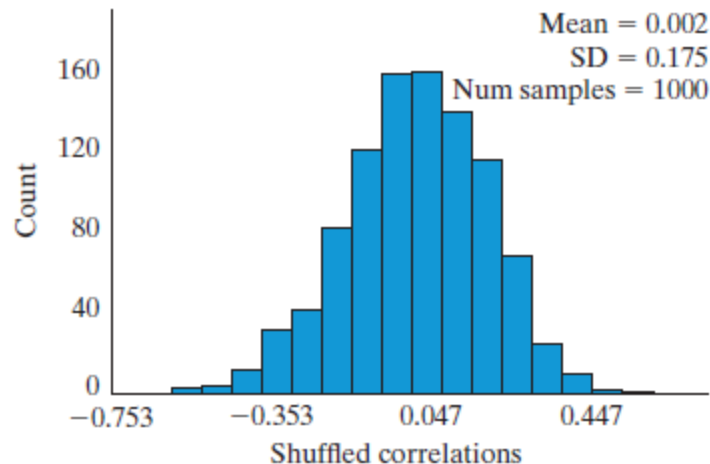
- Testing the slope of the regression line is equivalent to testing the correlation (same p-value, but obviously different confidence intervals since the statistics are different)
- Hence these hypotheses are equivalent.
  - $H_0: \beta = 0$     $H_a: \beta < 0$  (Slope)
  - $H_0: \rho = 0$     $H_a: \rho < 0$  (Correlation)
- Sample slope (b)   Population ( $\beta$ : beta)
- Sample correlation (r)   Population ( $\rho$ : rho)
- When we do the theory based test, we will be using the *t*-statistic which can be calculated from either the slope or correlation.

# Introduction

- Our null distributions are again bell-shaped and centered at 0 (for either correlation or slope as our statistic).

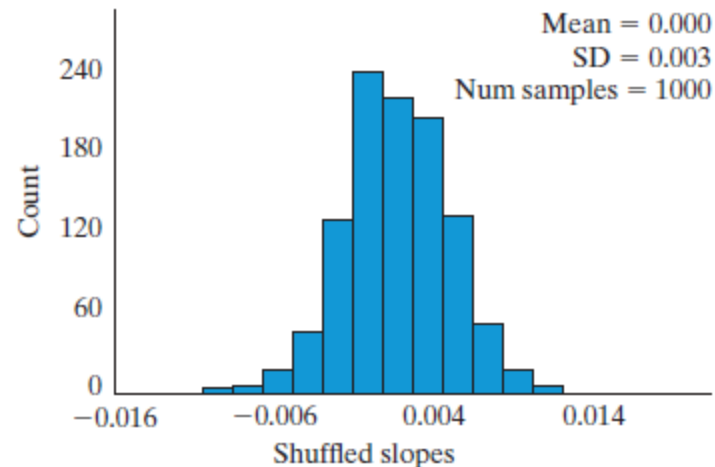
Example 10.2: Exercise and mood intensity

● Correlation ○ Slope ○ *t*-statistic



Example 10.4: GPA and nonacademic hours

○ Correlation ● Slope ○ *t*-statistic



The book on p549 finds a p value of 3.3% by simulation.

# Validity Conditions

- Under certain conditions, theory-based inference for correlation or slope of the regression line use  $t$ -distributions.
- We could use simulations or the theory-based methods for the slope of the regression line.
- We would get the same p-value if we used correlation as our statistic.

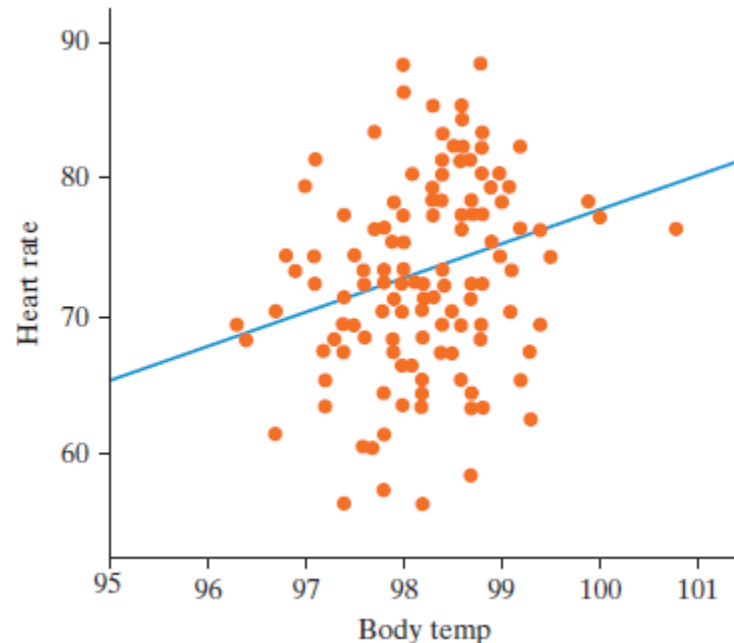
# Predicting Heart Rate from Body Temperature

*Example 10.5A*



# Heart Rate and Body Temp

- Earlier we looked at the relationship between heart rate and body temperature with 130 healthy adults
- Predicted Heart Rate =  $-166.3 + 2.44(\text{Temp})$
- $r = 0.257$

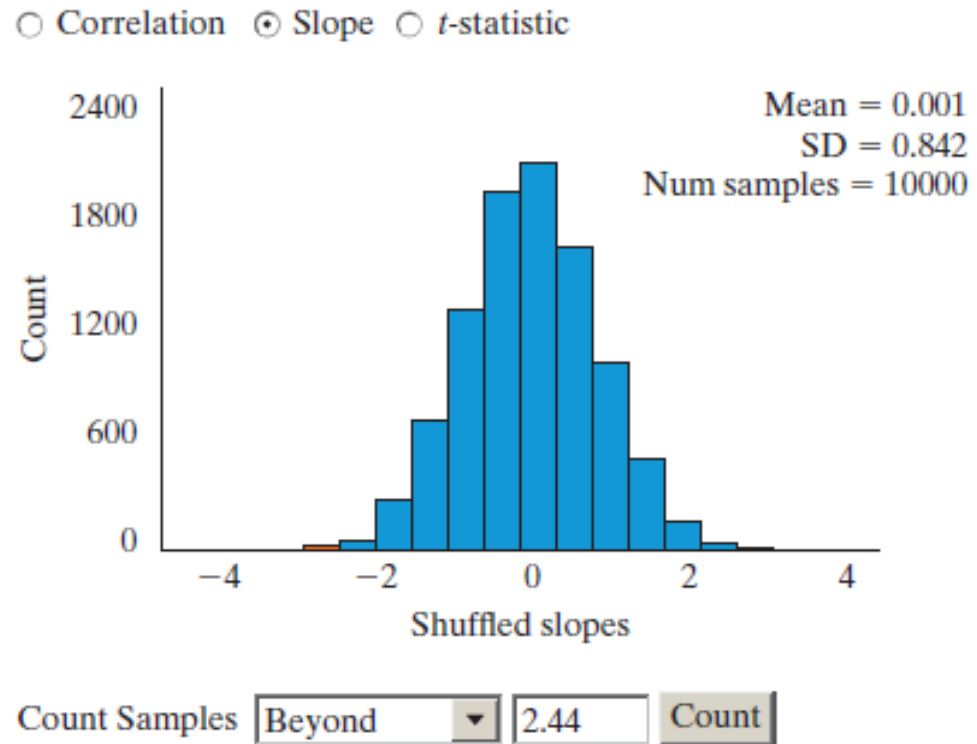


# Heart Rate and Body Temp

- We tested to see if we had convincing evidence that there is a positive association between heart rate and body temperature in the population using a simulation-based approach. (We will make it 2-sided this time.)
- **Null Hypothesis:** There is no association between heart rate and body temperature in the population.  $\beta = 0$
- **Alternative Hypothesis:** There is an association between heart rate and body temperature in the population.  $\beta \neq 0$

# Heart Rate and Body Temp

We get a very small p-value (0.0036). Anything as extreme as our observed slope of 2.44 happening by chance is very rare



# Heart Rate and Body Temp

- We can also approximate a 95% confidence interval

observed statistic  $\pm$  1.96 SD of statistic

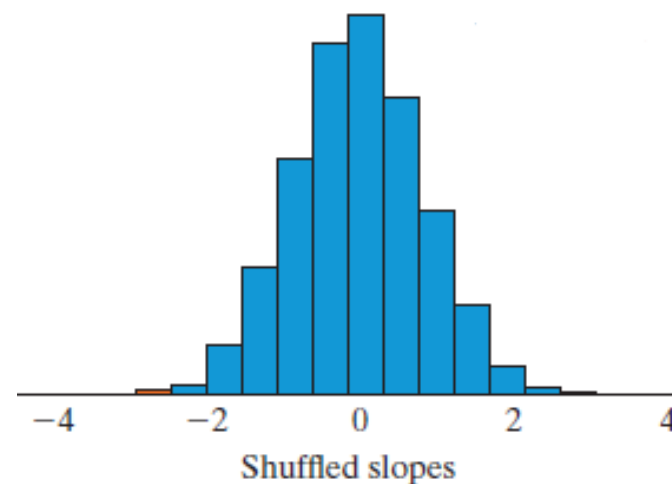
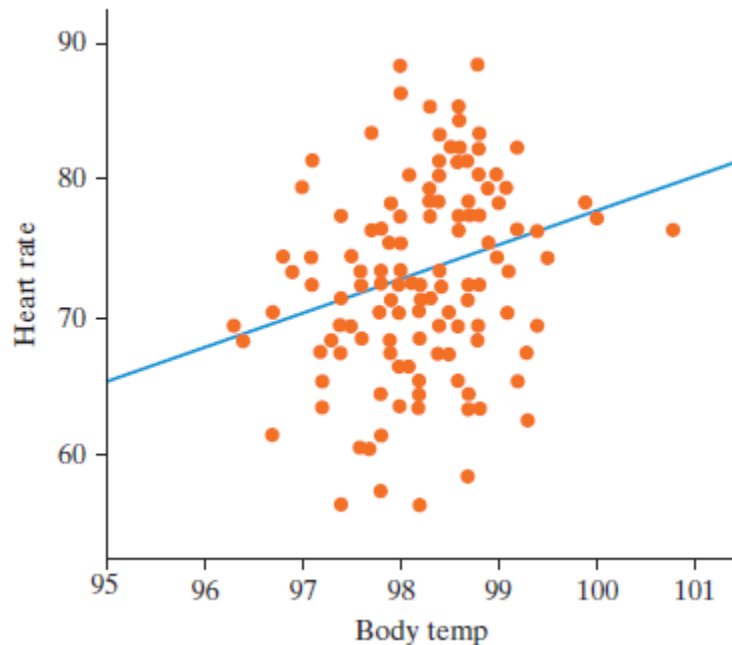
$$2.44 \pm 1.96(0.842) = 0.790 \text{ to } 4.09$$

- What does this mean?

We're 95% confident that, in the population of healthy adults, each 1° increase in body temp is associated with an increase in heart rate of between 0.790 to 4.09 beats per minute

# Heart Rate and Body Temp

- The theory-based approach should work well since the distribution has a nice bell shape
- Also check the scatterplot



# Heart Rate and Body Temp

- We will use the t-statistic to get our theory-based p-value.
- We will find a theory-based confidence interval for the slope.
- On p554, the book notes the formula
- $t = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$ .
- Here the t statistic is 2.97.
- The p-value is .36%. So the correlation is statistically significantly greater than zero.

# Smoking and Drinking

Example 10.5B

# Validity Conditions

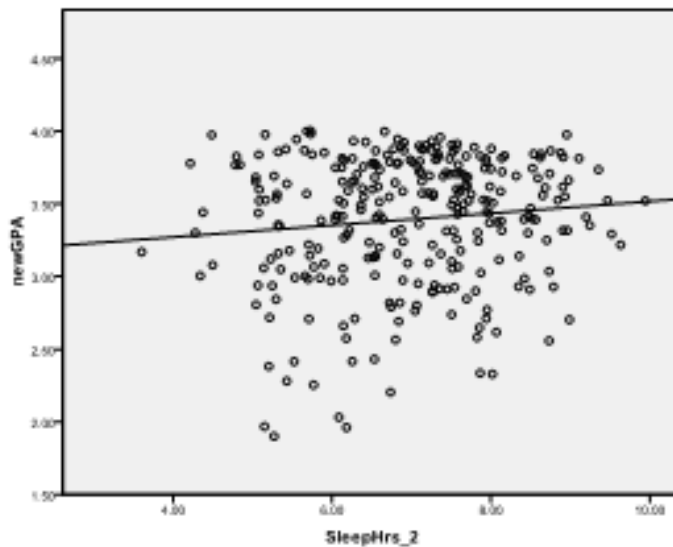
Remember our validity conditions for theory-based inference for slope of the regression equation.

1. The scatterplot should follow a linear trend.
2. There should be approximately the same number of points above and below the regression line (symmetry).
3. The variability of vertical slices of the points should be similar. This is called homoskedasticity.

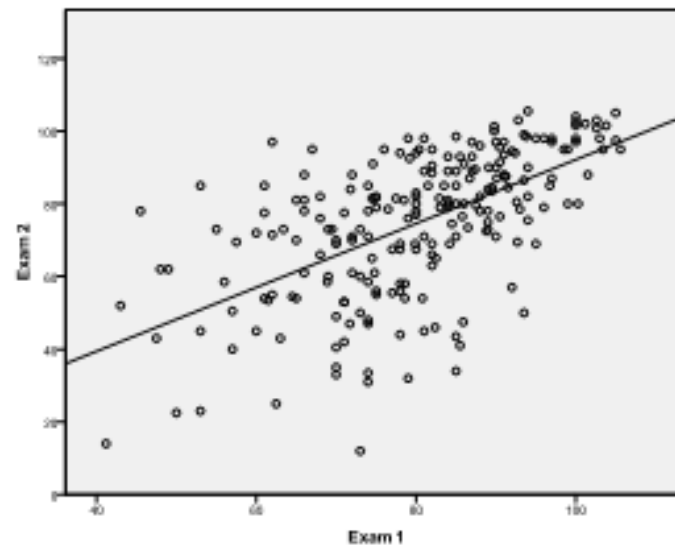


# Validity Conditions

- Let's look at some scatterplots that do not meet the requirements.



(a)



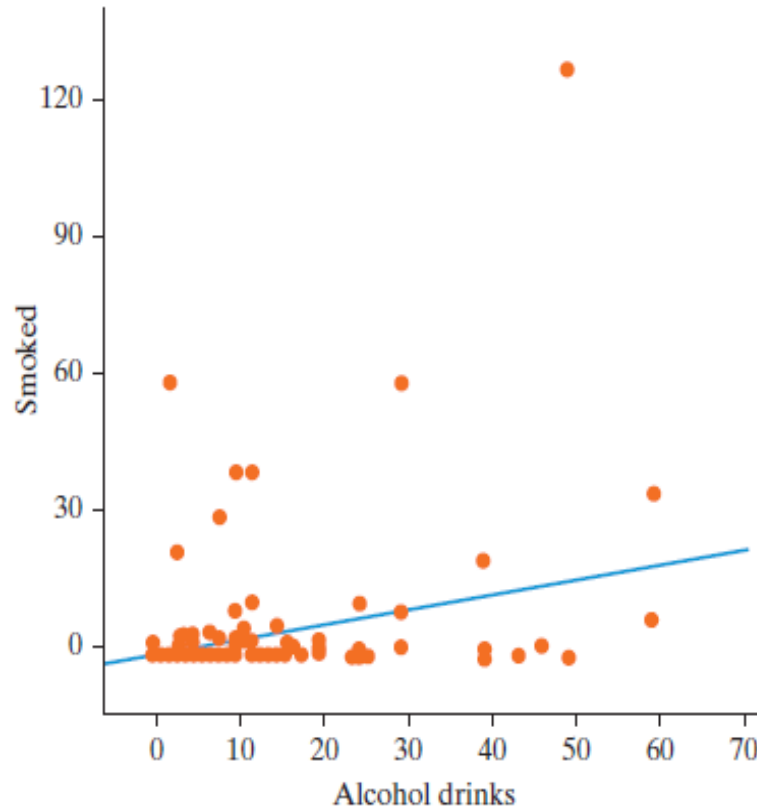
(b)

# Smoking and Drinking

The relationship between number of drinks and cigarettes per week for a random sample of students at Hope College.

The dot at (0,0)  
represents 524  
students

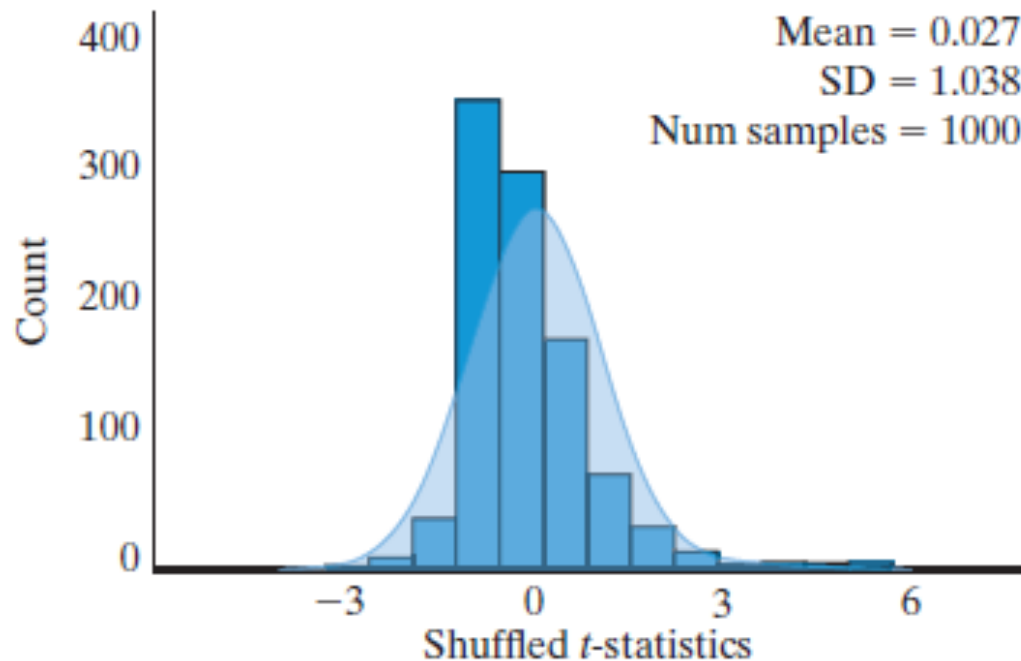
Are the conditions met?  
Hard to say. The book  
says no.



# Smoking and Drinking

- When the conditions are not met, applying simulation-based inference is preferable to theory-based t-tests and CIs.

○ Correlation ○ Slope ⊙ *t*-statistic



# Validity Conditions

- What do you do when validity conditions aren't met for theory-based inference?
  - Use the simulated-based approach.
- Another strategy is to “transform” the data on a different scale so conditions are met.
  - The logarithmic scale is common.
- One can also fit a different curve, not necessarily a line.