

# Stat 201: Statistics I

## Week 11



# Week 11

## Correlation and Regression

# Comparing samples from two populations

## Example

Consider the following data:

X	83	85	66	89	96	78
Y	90	90	66	86	99	85

How should this data be analyzed? It depends on the context of the data or what the data represents.

- Suppose the data is test scores from two different statistics classes.
  - One possible analysis would be a two-sample t-test comparing the mean test scores from each class.
- Suppose the data are scores from the midterm and the final for one set of students.
  - One possible analysis would be a matched pairs t-test comparing the mean difference between the midterm and the final for each student.

# Comparing samples from two populations, cont.

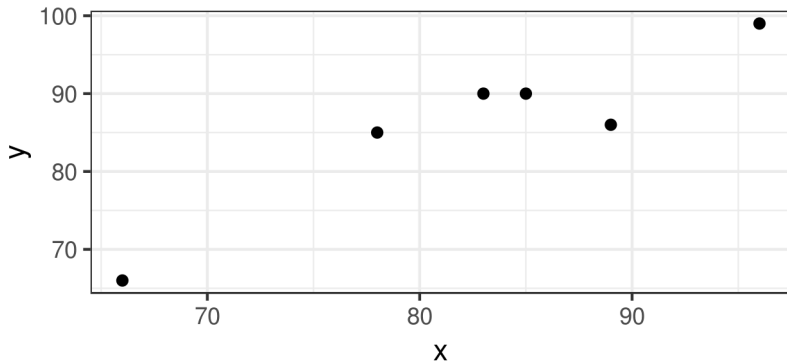
## Example

X	83	85	66	89	96	78
Y	90	90	66	86	99	85

- Suppose X is a student's score on the statistics final and Y is that student's yearly salary, in thousands of dollars, a year later.
  - It doesn't really make sense to compare means. The two samples represent entirely different kinds of data. There is no meaningful way to compare means.
  - It is useful to examine the association between the data. Is a higher test score associated with a higher salary? Or are the samples independent, values of one having no effect on values of the other?

# Scatterplots

A **scatterplot** is a useful way to visually examine paired data. Every data pair is represented by one point in the plot.



# Section 11.1

## Correlation

# Correlation

**Correlation** is the measure of how strongly associated values of paired data are with each other.

**Linear correlation** is the measure of how strong of a linear relationship (along a straight line) values of paired data have with each other.

- A positive correlation indicates that as one value in a pair increases the other will tend to increase.
- A negative correlation indicates that as one value in a pair increases the other will tend to decrease.
- No correlation indicates that the two values of a pair have no relationship with each other. They are independent.

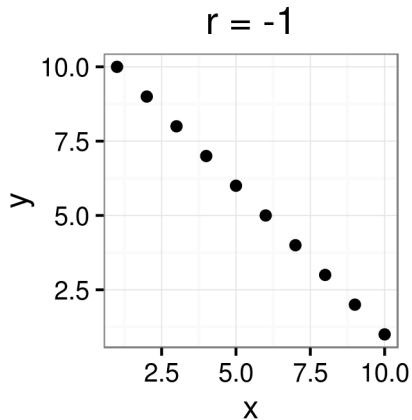
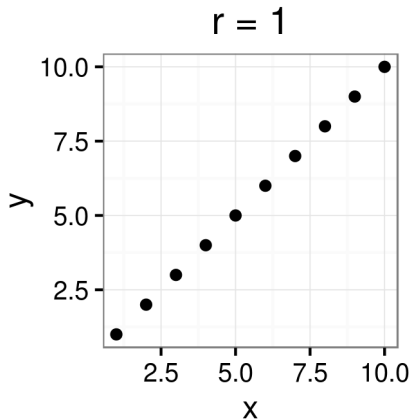
# Correlation coefficient

The sample **correlation coefficient**, denoted as  $r$ , is the numeric value for the strength of linear correlation between values of paired data.

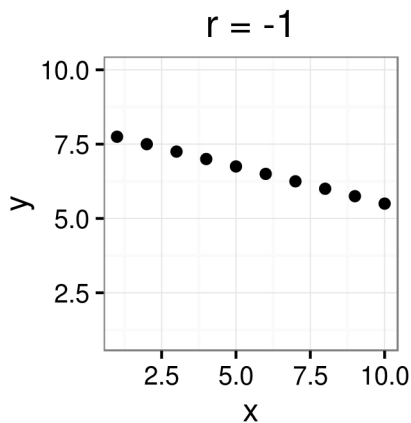
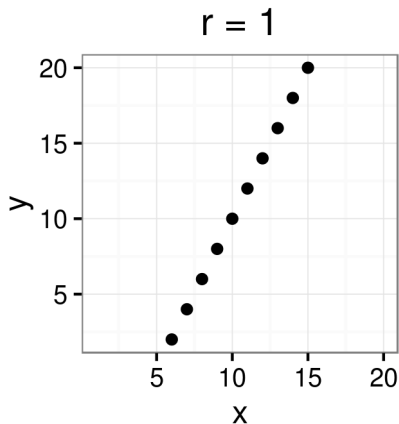
- $r$  is always between -1 and 1
- When  $r = 1$ , the samples have perfect positive correlation
- When  $r = -1$ , the samples have perfect negative correlation
- When  $r = 0$ , the samples are perfectly independent
- Most samples will have some other value for  $r$ .
- The order of the values, i.e.  $(x, y)$  vs.  $(y, x)$ , have no effect of the value of  $r$ .
- The units of the values also have no effect on  $r$ . A correlation on height will be the same whether it is measured in inches, feet or meters, for example.



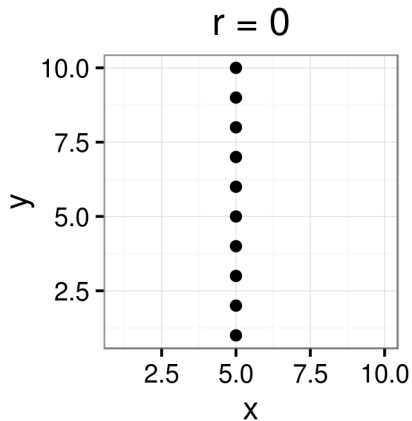
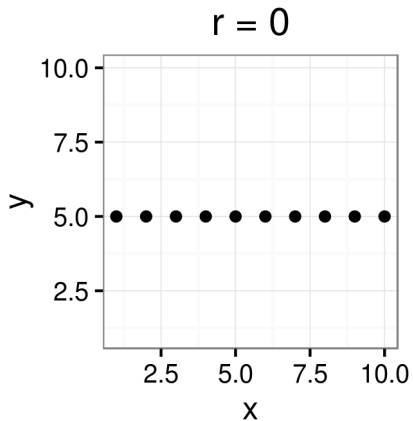
# Perfect correlation



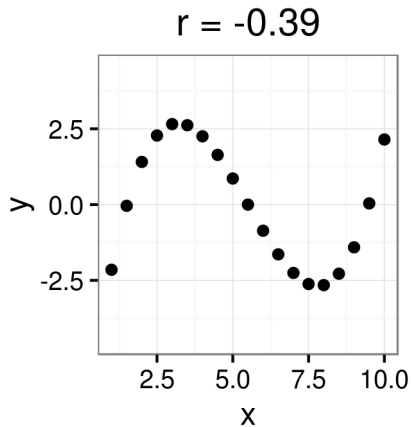
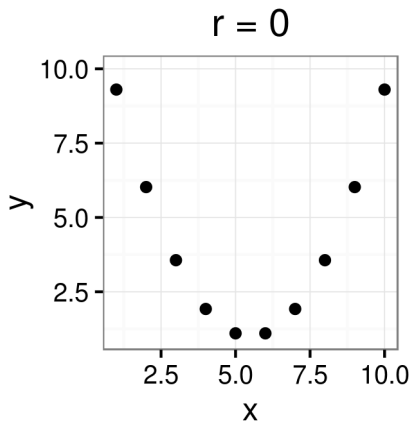
# Perfect correlation, cont.



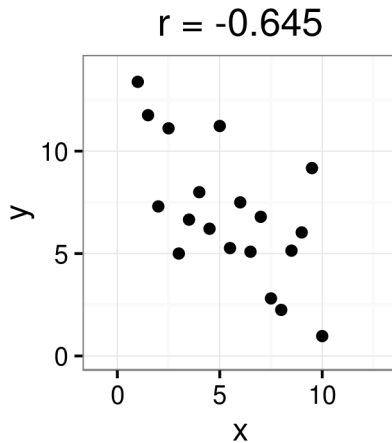
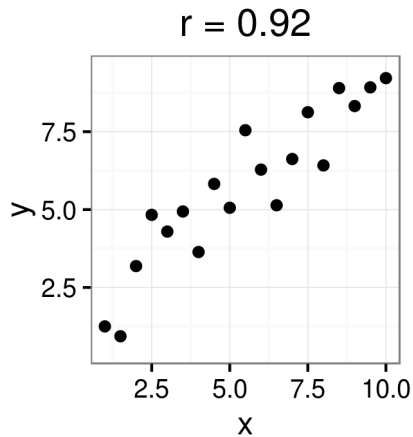
# Zero correlation



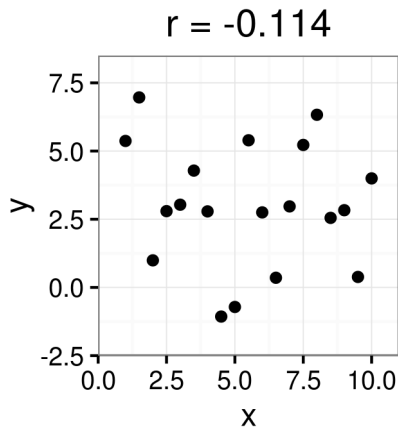
# Non-linear correlation



# “Real world” correlation

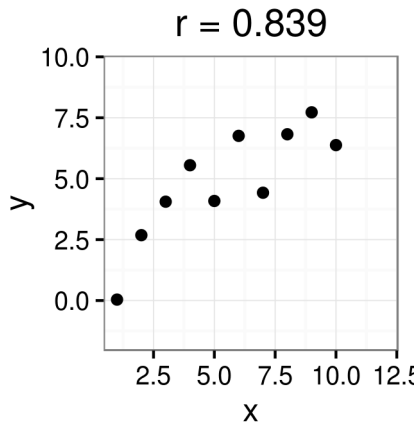
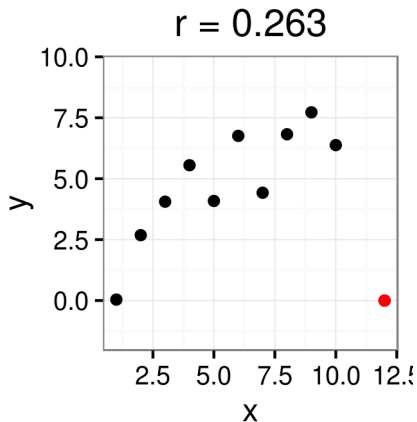


# “Real world” correlation, cont.



# Correlation and outliers

The correlation coefficient  $r$  is sensitive to outliers (very unusual values).



# Calculating correlation coefficient

To calculate  $r$  one could use the following formula:

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{n(\sum x^2) - (\sum x)^2} \sqrt{n(\sum y^2) - (\sum y)^2}}$$

- While this formula will give you a value for  $r$  (after much work), it doesn't help much in understanding how the correlation coefficient works.
- Since technology should be used to calculate the correlation calculation anyway, it would be useful to look at  $r$  defined with a more informative formula.



# Calculating correlation coefficient, cont.

Consider: ( $\propto$  means “is proportional to”)

$$r \propto \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- If when  $x$  values are big (greater than  $\bar{x}$ ) then  $y$  values are big (greater than  $\bar{y}$ ), and vice versa ( $x < \bar{x} \rightarrow y < \bar{y}$ ), then  $r$  is proportional to the sum of positive numbers, i.e.  $r > 0$ .
- If when  $x$  values are big (greater than  $\bar{x}$ ) then  $y$  values are small (less than  $\bar{y}$ ), and vice versa ( $x < \bar{x} \rightarrow y > \bar{y}$ ), then  $r$  is proportional to the sum of negative numbers, i.e.  $r < 0$ .
- If when  $x$  values are far away from  $\bar{x}$  then  $y$  values are also far away from  $\bar{y}$  (whether positive or negative), then  $r$  is proportional to the sum of larger numbers and will be closer to 1 or -1.

# Correlation hypothesis tests

The sample correlation coefficient  $r$  is an unbiased estimator for the population correlation coefficient  $\rho$  (rho).

Like population proportions and means, hypothesis tests can be conducted of the parameter  $\rho$  to determine if a population of paired data is correlated.

- $H_0 : \rho = 0$

- $H_a : \rho \neq 0$

(It is uncommon to use the other forms of the alternative hypothesis)

- Test statistic:  $t = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}} \quad df = n - 2$

- If  $p < \alpha$ , reject null hypothesis. There is evidence that the samples are correlated.

# Strength of correlation

Even with statistically significant correlation, it might be useful to specify the strength of the correlation. While, like many things in statistics, the description of strength can be subjective and dependent of the context of the data, the following rule of thumb may be used:

- $|r| \geq 0.7$  is a strong correlation
- $0.4 \leq |r| < 0.7$  is a moderate correlation
- $|r| < 0.4$  is a weak correlation

Thus,  $r = 0.5$  indicates a moderate positive correlation and  $r = -0.8$  indicates a strong negative correlation.

# Correlation, example

## Example

In 1886, Sir Francis Galton, a British sociologist, published the paper “Regression towards Mediocrity in Hereditary Statures”, in which he examined the heights of parents and their adult children. The core of modern uses of correlation and regression come from this paper (he also invented standard deviation). His data for fathers and sons is in the file “Galton-father-son.csv” on D2L.

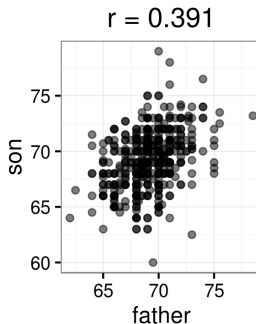
Is there correlation between heights of fathers and their adult sons? Test at  $\alpha = 0.05$  level of significance.

- $H_0 : \rho = 0$   
 $H_a : \rho \neq 0$
- $r = 0.39131736$   
 $p \ll 0.0001$

# Correlation, example

## Example

- $p < 0.0001 < \alpha = 0.5$ . Reject null hypothesis
- There is evidence that there is a correlation between the heights of fathers and their adult sons. However,  $r < 0.4$  indicates weak correlation.



# Coefficient of determination

The **coefficient of determination**, designated by  $R^2$ , represents the proportion of variation in one variable that is explained by the association with the other variable.

- The coefficient of determination is the correlation coefficient squared,  $R^2 = r^2$

## Example

In the previous example, heights of fathers and sons had a correlation coefficient of  $r = 0.391$ . Thus,  $R^2 = (0.391)^2 = 0.153$ .

About 15% of the variation of the heights of adult men can be explained by the association with their fathers heights.

# Cautions

If a sample of paired data has a correlation coefficient that is zero or very low, that does not necessarily mean that there is not a association between the variables, only that there is not a *linear* association.

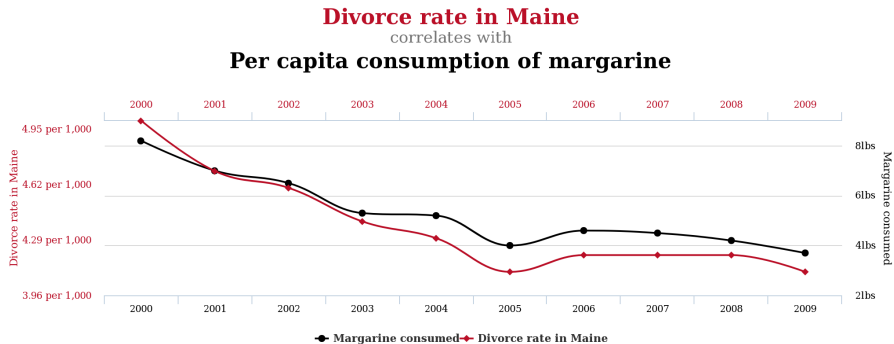
!!!

Correlation does not imply causation.

# Spurious correlations

From the *Spurious Correlations* website:

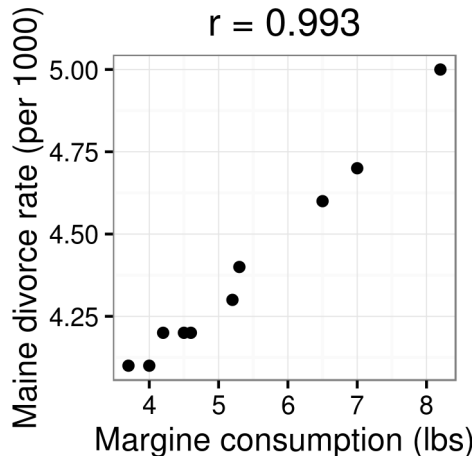
<http://www.tylervigen.com/spurious-correlations>



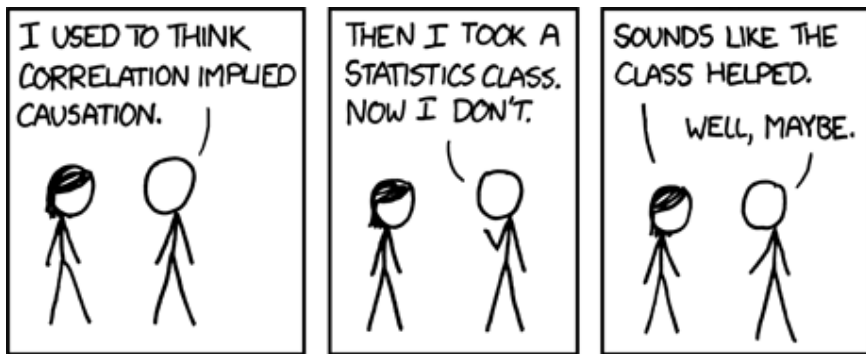
tylervigen.com



# Spurious correlations, cont.



# Correlation does not imply causation



# Section 11.2

## Regression

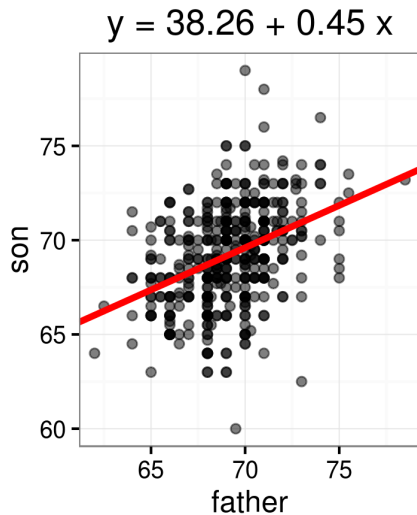
# Regression

Once it has been determined that two variables have a linear relationship, that is they have an association described by a straight line, the next logical question is what is that line.

**Regression** is the statistical technique for finding the line that best describes a linear relationship between two paired variables.

The line found is known as the **regression line** or the line of best fit.

# Regression, example



# Algebra review, lines

The equation for a line generally has the following form:

$$y = b + mx$$

- $b$  is the  $y$ -intercept, or where the line crosses the  $y$ -axis ( $x = 0$ ).
- $m$  is the slope of the line. It is the amount the  $y$  value increases as the  $x$  value increase by one.

# Regression population models

A linear relationship between populations of variables  $X$  and  $Y$  can be described by:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- This equation is the **model**, a possible structure for the relationship between the data.
- $X$  is known as the predictor variable, or the explanatory variable, or the independent variable.
- $Y$  is the response variable, or the dependent variable.
- $\beta_0$  and  $\beta_1$  are the intercept and slope of a line describing the association of  $X$  and  $Y$ .
- Like other population parameters,  $\beta_0$  and  $\beta_1$  are thought of as fixed, but unknown
- $\epsilon$  (epsilon) is a random error term. It is usually assumed that  $\epsilon \sim N(0, \sigma^2)$ .

# Regression lines

The regression line for a sample paired data  $(x, y)$  describes the linear relationship between  $x$  and  $y$ . It is given by the equation:

$$\hat{y} = b_0 + b_1x$$

- $\hat{y}$  is the estimated  $y$  value for a given  $x$ , or the predicted response.
- $b_1$ , the slope of the regression line, is the amount of change in the predicted response for each unit increase in  $x$
- $b_1$  can be calculated from the correlation coefficient,  $b_1 = r \frac{s_y}{s_x}$
- $b_0$ , the  $y$ -intercept of the regression line, is usually not of much interest.



# Hypothesis tests for regression

A simple linear regression model has two population parameters,  $\beta_0$  and  $\beta_1$ , that are estimated by sample statistics,  $b_0$  and  $b_1$ . As with other estimated parameters, an hypothesis test can be conducted to understand the population parameters. Generally, only the slope ( $\beta_1$ ) is tested.

- $H_0 : \beta_1 = 0$
- $H_a : \beta_1 \neq 0$   
(this is most common, though the other forms can be used.)
- $b_1$  follows a  $t$  distribution with  $n - 1$  degrees of freedom.
- This is the exact same test as the correlation test.

# Predictions

If there is a valid regression equation ( $\beta_1 \neq 0$ ), it can be used to make predictions of the response variable for given values of the predictor variable. Replace  $x$  in the equation with the given predictor value and calculate the predicted response  $\hat{y}$ .

However, if there is no correlation between values of paired data (if failed to reject null hypothesis for correlation test), then the best predictor for the response variable is simply  $\bar{y}$ .

# Predictions, example

## Example

The regression line equation from the Galton data, for fathers height as predictor  $x$  and sons height as response  $y$ , is

$$\hat{y} = 38.26 + 0.45x$$

What is the predicted adult height of a son whose father is 68 inches tall?

- $\hat{y} = 38.26 + 0.45 \times 68 = 68.86$  inches

# Extrapolating

Predictions for predictor values outside the range of  $x$  values used to find the regression line are highly suspect. This is known as **extrapolating** and should be avoided.

## Example

For the 2016-2017 season, NBA teams had field goal percentages between 43.5% and 49.5%. A regression line of the relationship between FG% ( $x$ ) and numbers of games won ( $y$ ) is

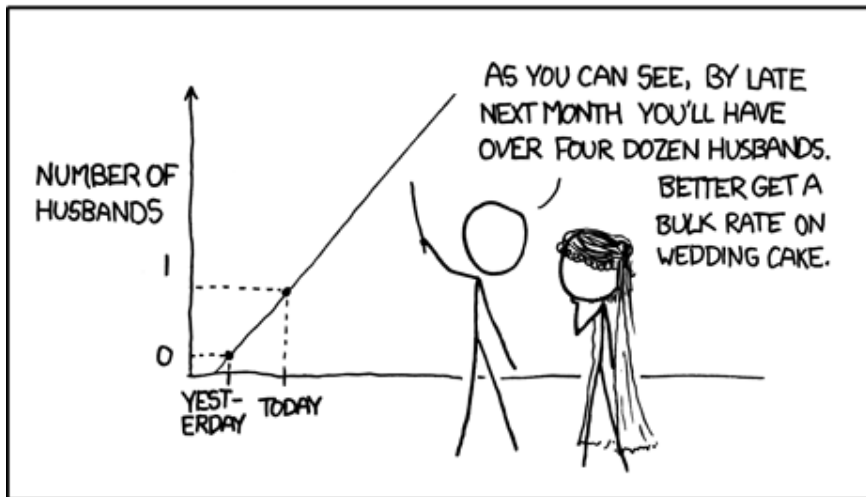
$$\hat{y} = -208.69 + 5.46x$$

The GM for the Timberwolves wants to know how many games the team would win if they could get their FG% up to 60%, based on this data.

- $\hat{y} = -208.69 + 5.46 \times 60 = 118.91$  games won

# Extrapolating, example

## MY HOBBY: EXTRAPOLATING



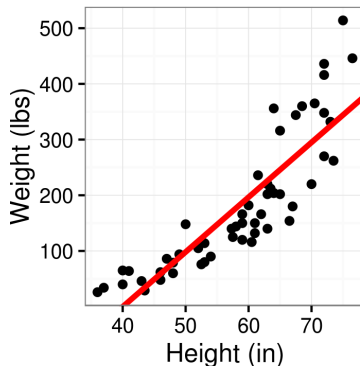
# Regression, example

## Example

The Department of Natural Resources wishes to track the weight of bears in the wild. While it is very difficult to weigh a bear, it is fairly easy to estimate the length of bear using photos. The data set “bears.csv” on D2L contains measurements made from anesthetized wild bears.

- Find a relationship, if any, between the length and weight of bears using the data.
- What is the best predicted weight of a bear thought to be 71 inches long?
- Would it be appropriate to predict the weight of a bear 39 inches long? 89 inches?

# Regression, example



## Example

- $R^2 = 0.747$ : About 75% of the variation in bear weight is explained by the association with bear height.

# Regression, example

## Example

- Regression line equation:  $\hat{y} = -393.84 + 9.84x$
- For every inch long a bear is, its weight will increase by about 10 lbs.
- A bear that is 71 inches long should weigh:

$$\hat{y} = -393.84 + 9.84 \times 71 = 304.7282 \text{ lbs.}$$

- The  $x$  values (bear length) used in the regression model are in the interval (36, 76.5).
- Thus, it would be appropriate to predict the weight of a bear that is 39 inches long, but not a bear that is 89 inches long.