

# Examining Numerical Data

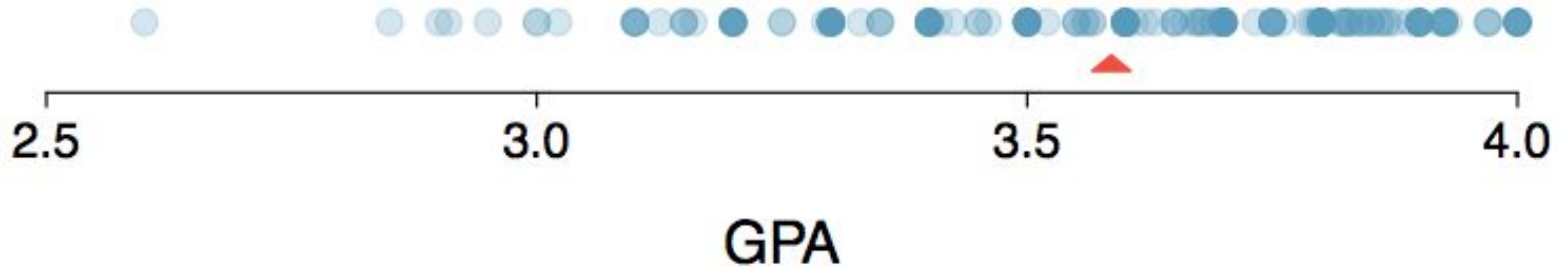
# Dot Plots

Useful for visualizing one numerical variable. Darker colors represent areas where there are more observations.



How would you describe the distribution of GPAs in this data set? Make sure to say something about the center, shape, and spread of the distribution.

# Dot Plots & Mean



The *mean*, also called the *average* (marked with a triangle in the above plot), is one way to measure the center of a *distribution* of data.

The mean GPA is 3.59.

# Mean

The *sample mean*, denoted as  $\bar{x}$ , can be calculated as

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n},$$

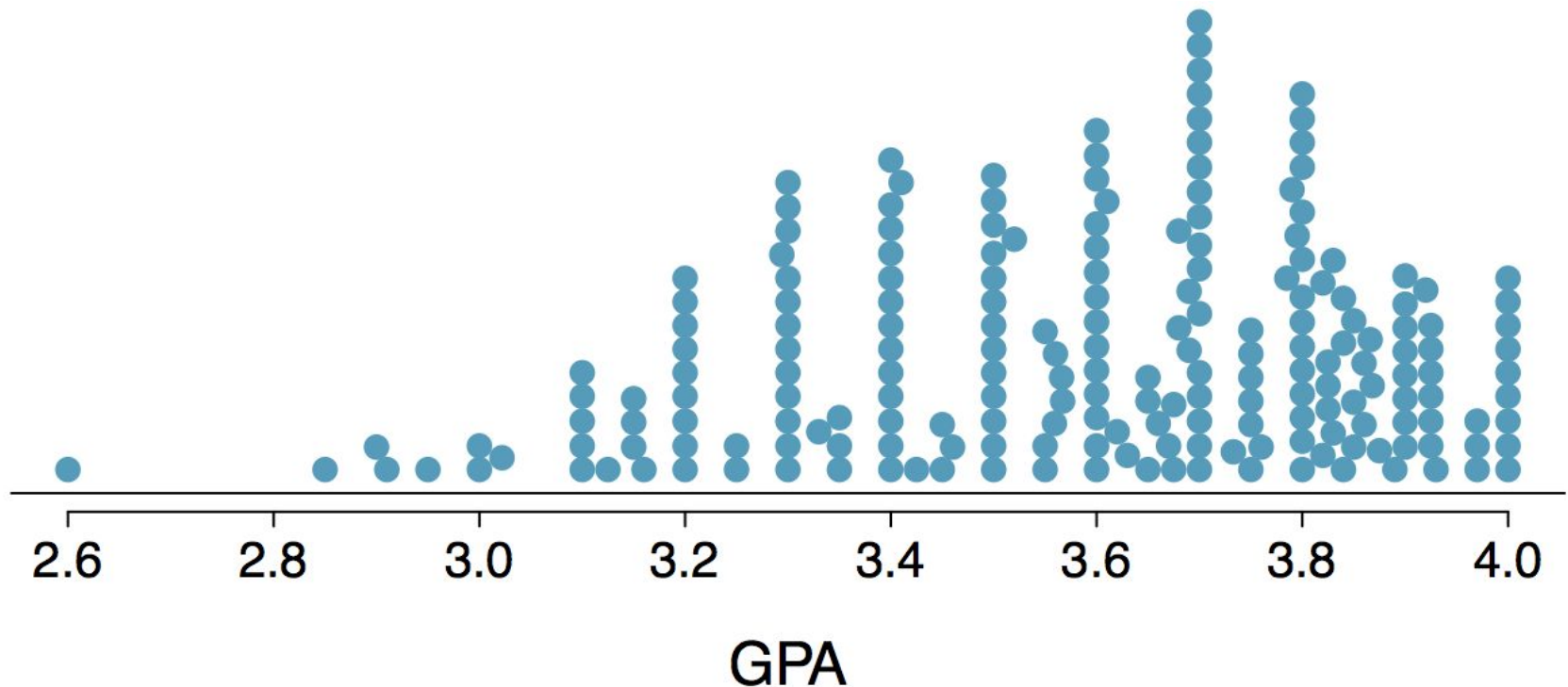
where  $x_1, x_2, \dots, x_n$  represent the  $n$  observed values.

The sample mean is a *sample statistic*, and serves as a *point estimate* of the population mean. This estimate may not be perfect, but if the sample is good (representative of the population), it is usually a pretty good estimate.

***What does it mean to be a “pretty good estimate”?***

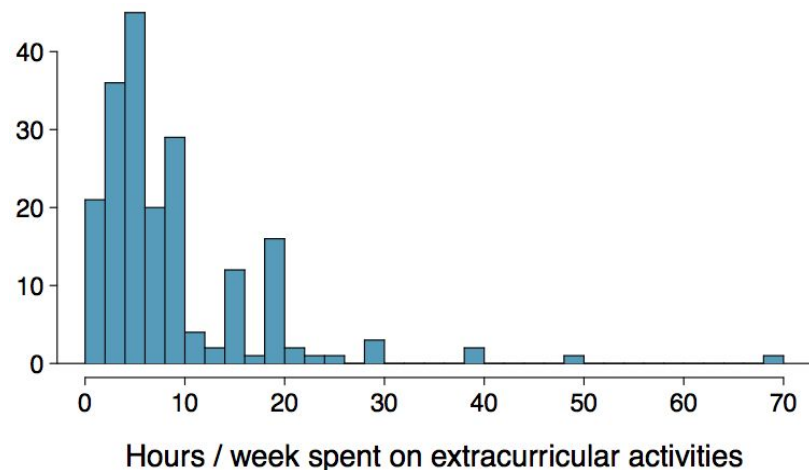
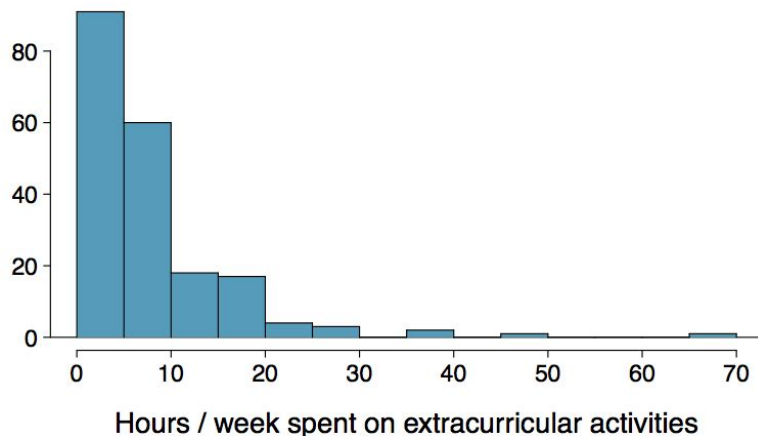
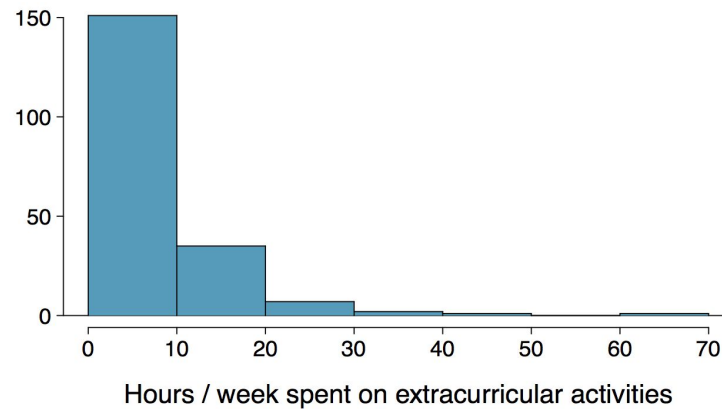
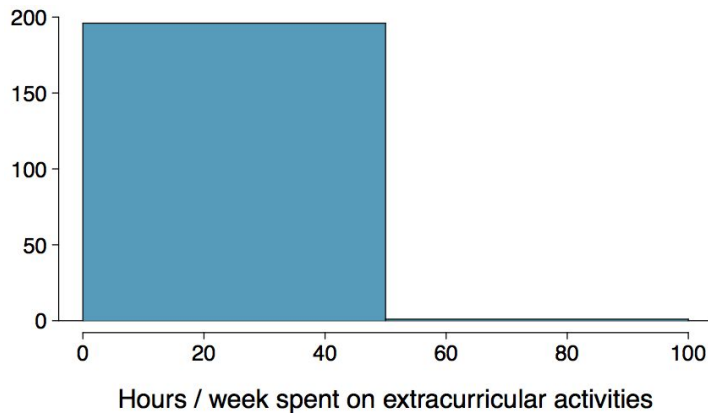
# Stacked Dot Plot

Higher bars represent areas where there are more observations, makes it a little easier to judge the center and the shape of the distribution.



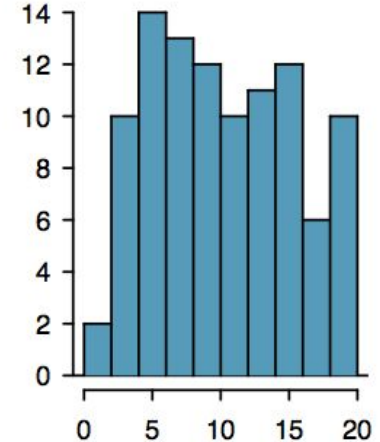
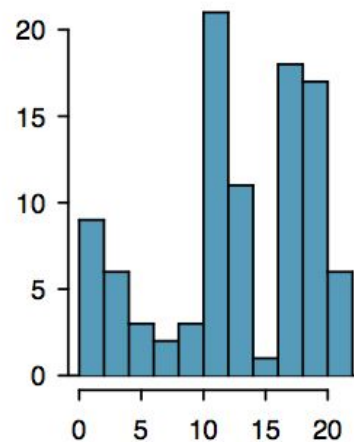
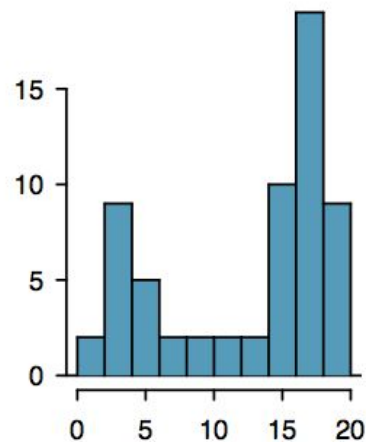
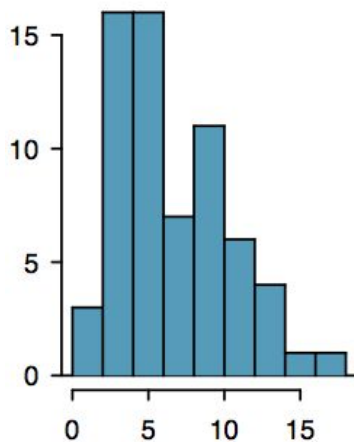
# Bin Width

Which one(s) of these histograms are useful? Which reveal too much about the data? Which hide too much?



# Shape of a Distribution: Modality

Does the histogram have a single prominent peak (*unimodal*), several prominent peaks (*bimodal/multimodal*), or no apparent peaks (*uniform*)?

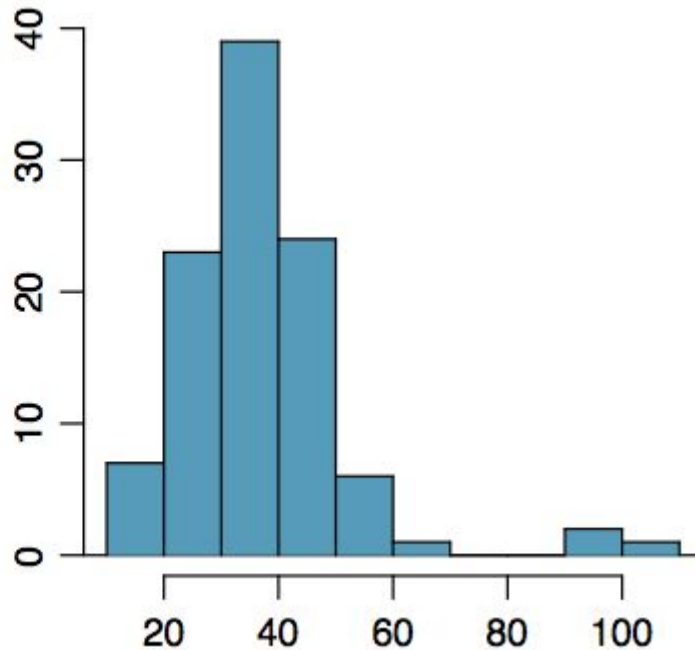
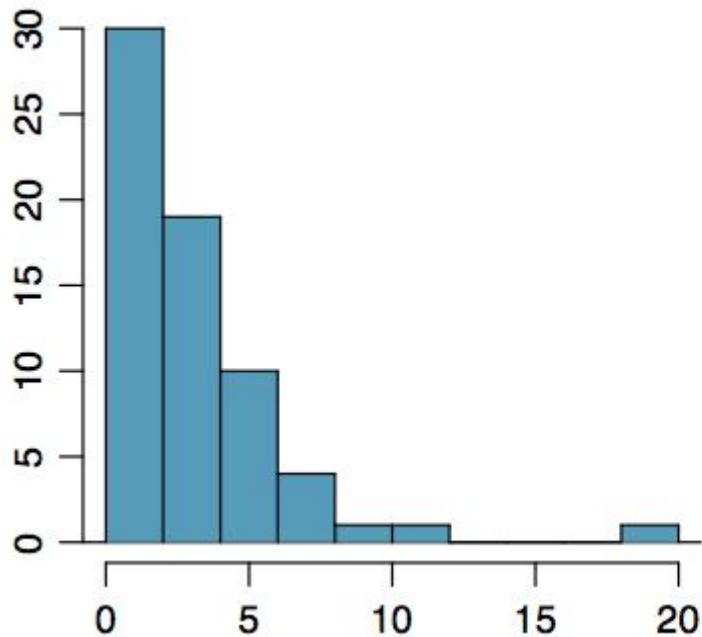


---

**Note:** In order to determine modality, step back and imagine a smooth curve over the histogram -- imagine that the bars are wooden blocks and you drop a limp spaghetti over them, the shape the spaghetti would take could be viewed as a smooth curve.

# Shape of a Distribution: Unusual Observations

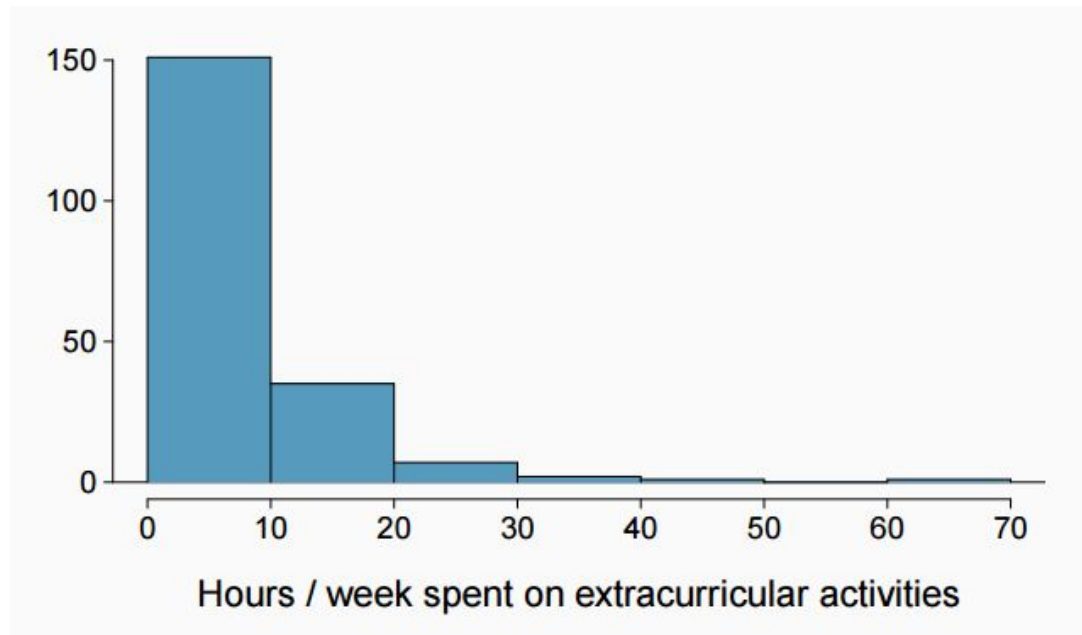
Are there any unusual observations or potential *outliers*?





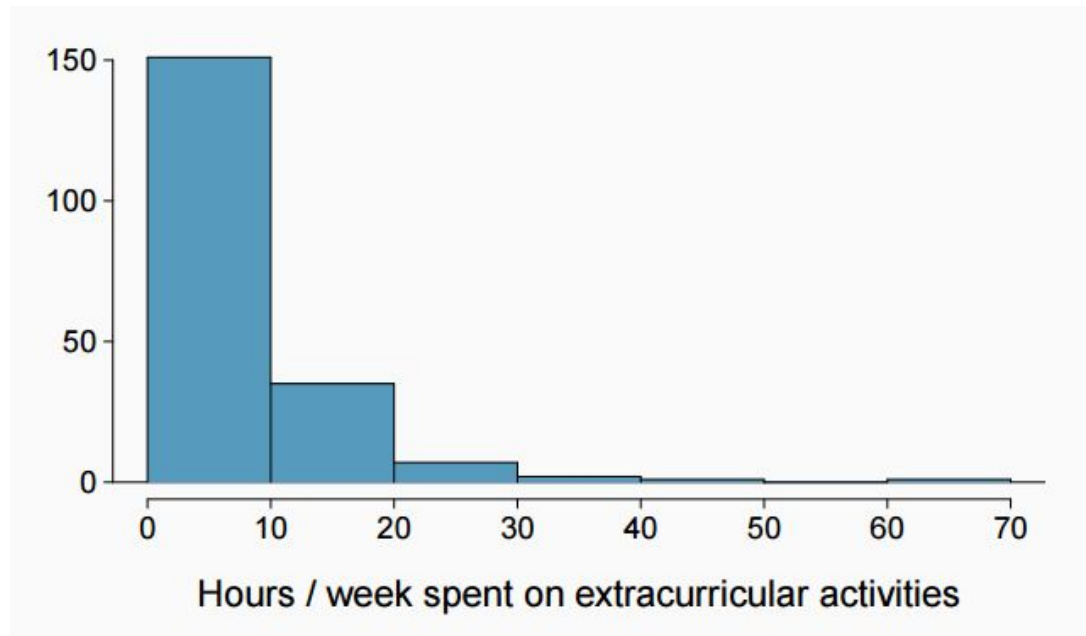
# Extracurricular activities

How would you describe the shape of the distribution of hours per week students spend on extracurricular activities?



# Extracurricular activities

How would you describe the shape of the distribution of hours per week students spend on extracurricular activities?



*Unimodal and right skewed, with a potentially unusual observation at 60 hours/week.*

# Commonly observed shapes of distributions

## Modality

unimodal



bimodal



multimodal



uniform



## Skewness

right skew



left skew



symmetric



# Variance

Variance is roughly the average squared deviation from the mean.

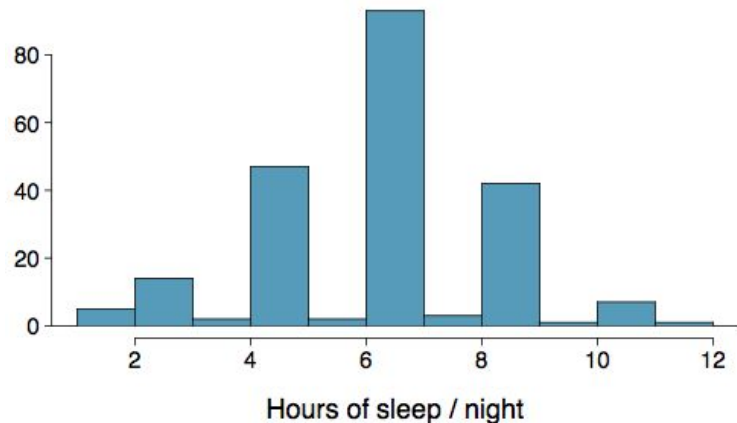
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

# Variance

Variance is roughly the average squared deviation from the mean.

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

- The sample mean is  $\bar{x} = 6.71$ , and the sample size is  $n = 217$ .

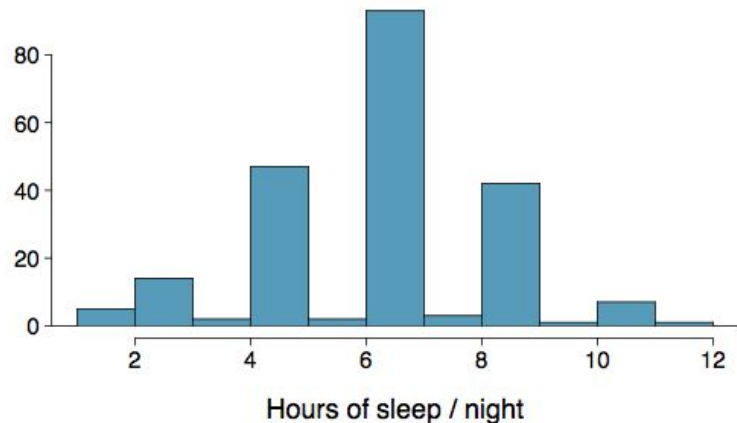


# Variance

**Variance** is roughly the average squared deviation from the mean.

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

- The sample mean is  $\bar{x} = 6.71$ , and the sample size is  $n = 217$ .
- The variance of amount of sleep students get per night can be calculated as:



$$s^2 = \frac{(5 - 6.71)^2 + (9 - 6.71)^2 + \dots + (7 - 6.71)^2}{217 - 1} = 4.11 \text{ hours}^2$$

# Variance (cont.)

Why do we use the squared deviation in the calculation of variance?

# Variance (cont.)

Why do we use the squared deviation in the calculation of variance?

- To get rid of negatives so that observations equally distant from the mean are weighed equally.
- To weigh larger deviations more heavily.
- The Euclidean distance is appropriate for work in higher dimensions; remember your linear algebra



# Standard Deviation

The *standard deviation* is the square root of the variance, and has the same units as the data.

$$s = \sqrt{s^2}$$

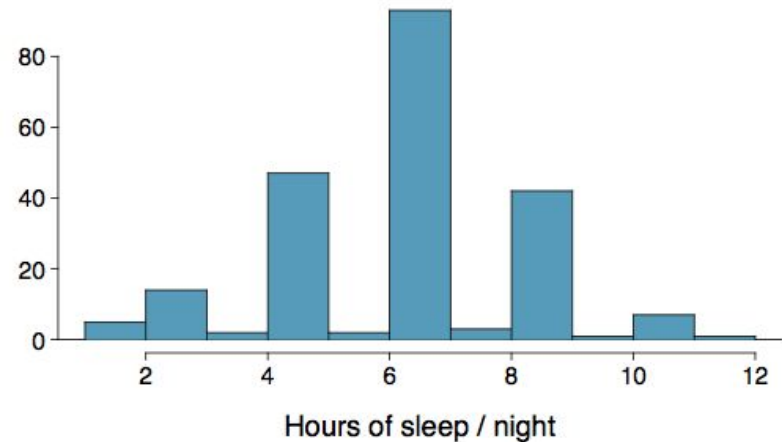
# Standard Deviation

The *standard deviation* is the square root of the variance, and has the same units as the data.

$$s = \sqrt{s^2}$$

- The standard deviation of amount of sleep students get per night can be calculated as:

$$s = \sqrt{4.11} = 2.03 \text{ hours}$$



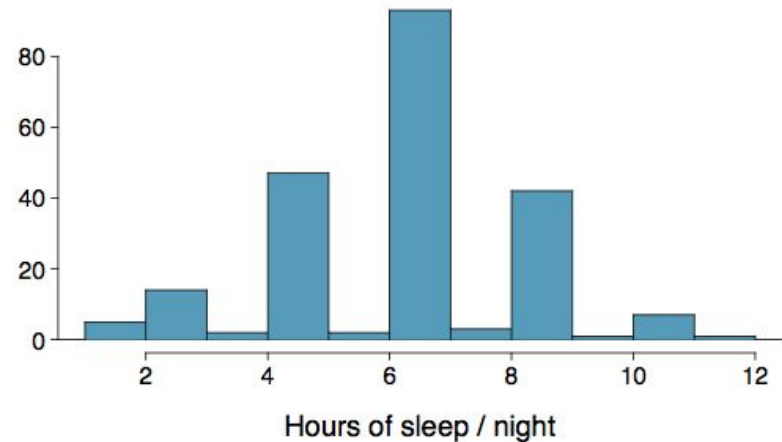
# Standard Deviation

The *standard deviation* is the square root of the variance, and has the same units as the data.

$$s = \sqrt{s^2}$$

- The standard deviation of amount of sleep students get per night can be calculated as:

$$s = \sqrt{4.11} = 2.03 \text{ hours}$$



- We can see that all of the data are within 3 standard deviations of the mean.

# Median

The *median* is the value that splits the data in half when ordered in ascending order.

0, 1, 2, 3, 4

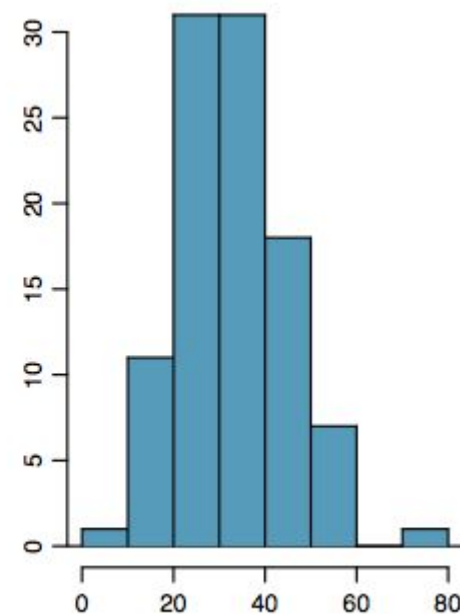
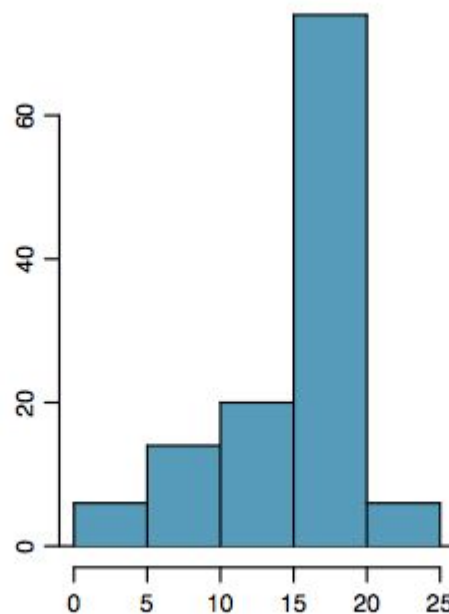
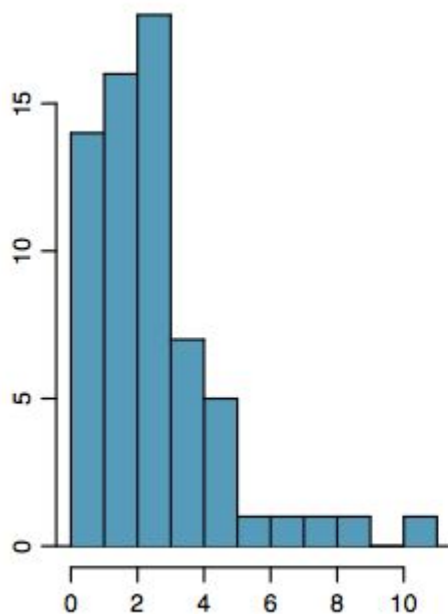
If there are an even number of observations, then the median is the average of the two values in the middle.

$$0, 1, \underline{2}, \underline{3}, 4, 5 \rightarrow \frac{2 + 3}{2} = 2.5$$

Since the median is the midpoint of the data, 50% of the values are below it. Hence, it is also the *50th percentile*.

# Shape of a Distribution: Skewness

Is the histogram *right skewed*, *left skewed*, or *symmetric*?



---

**Note:** Histograms are said to be skewed to the side of the long tail.

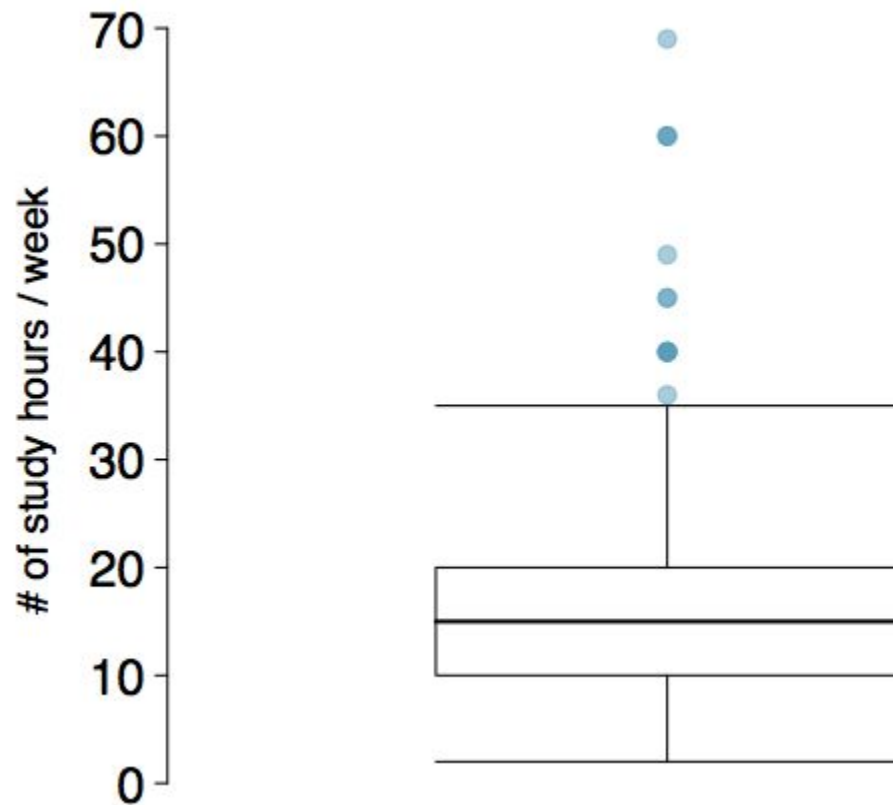
# Q1, Q3, and IQR

- The 25th percentile is also called the first quartile, *Q1*.
- The 50th percentile is also called the median.
- The 75th percentile is also called the third quartile, *Q3*.
- Between Q1 and Q3 is the middle 50% of the data. The range these data span is called the *interquartile range*, or the *IQR*.

$$IQR = Q3 - Q1$$

# Box Plot

The box in a *box plot* represents the middle 50% of the data, and the thick line in the box is the median.



# Whiskers and Outliers

*Whiskers* of a box plot can extend up to  $1.5 \times \text{IQR}$  away from the quartiles.

$$\text{max upper whisker reach} = Q3 + 1.5 \times \text{IQR}$$

$$\text{max lower whisker reach} = Q1 - 1.5 \times \text{IQR}$$

$$\text{IQR: } 20 - 10 = 10$$

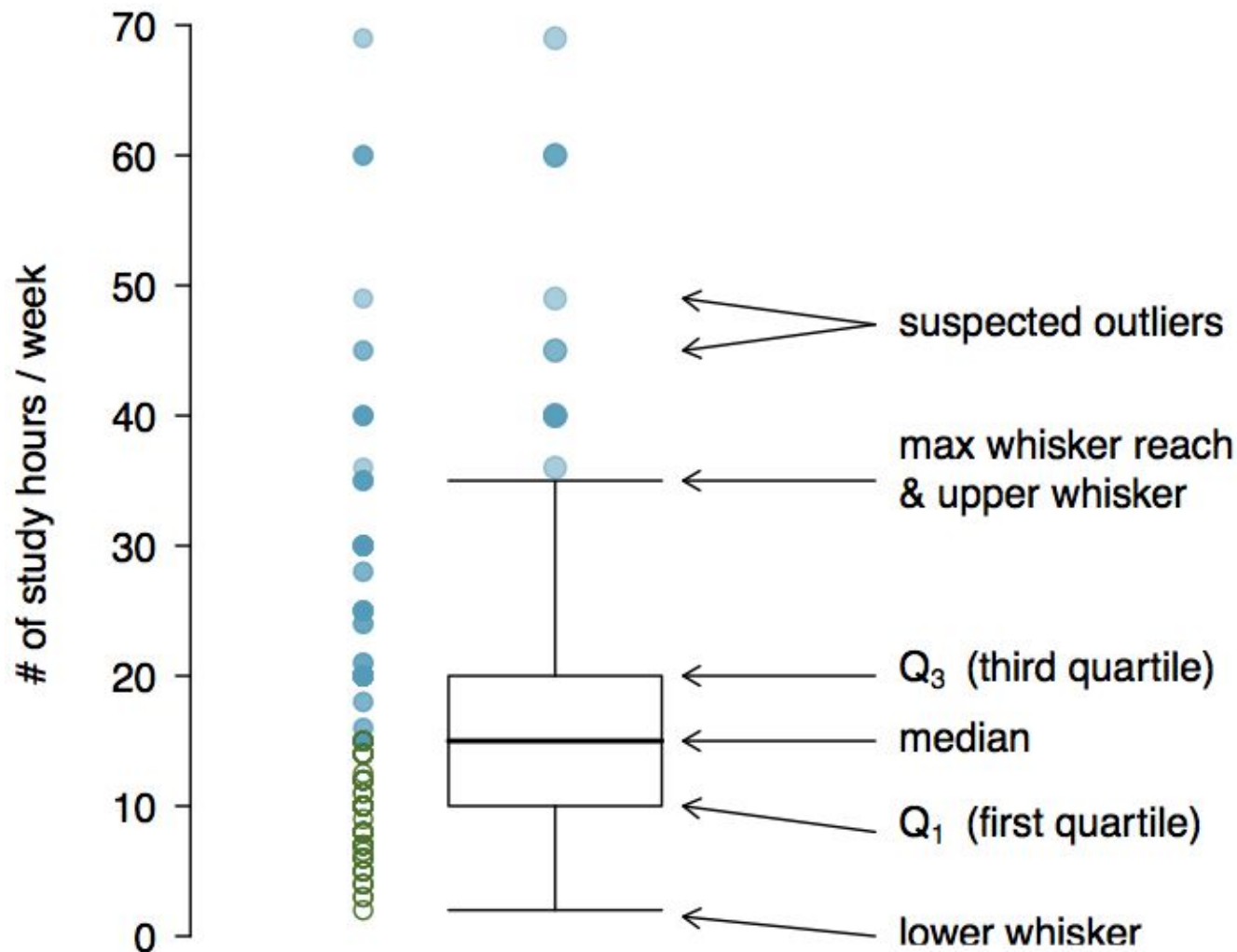
$$\text{max upper whisker reach} = 20 + 1.5 \times 10 = 35$$

$$\text{max lower whisker reach} = 10 - 1.5 \times 10 = -5$$

A potential *outlier* is defined as an observation beyond the maximum reach of the whiskers. It is an observation that appears extreme relative to the rest of the data.



# Anatomy of a Box Plot



# More on Outliers

Why is it important to look for outliers?

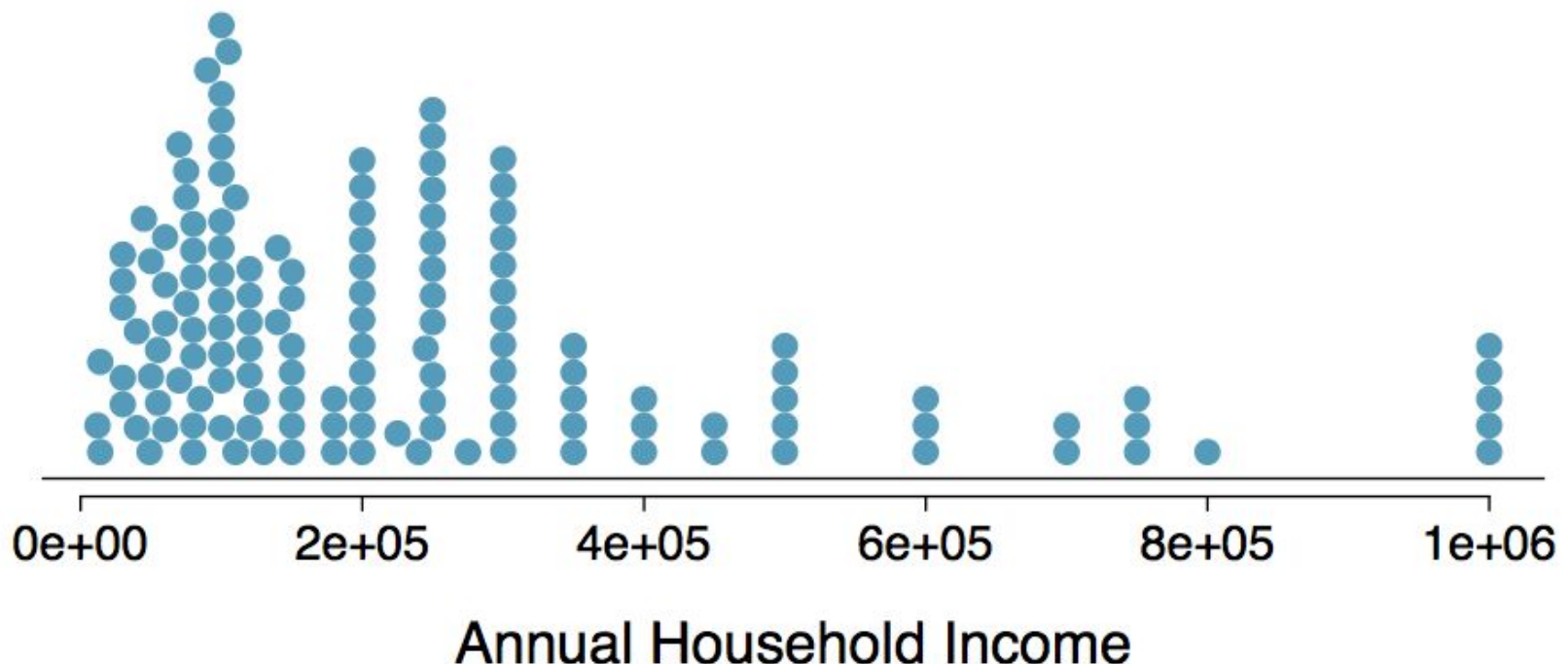
# Outliers (cont.)

Why is it important to look for outliers?

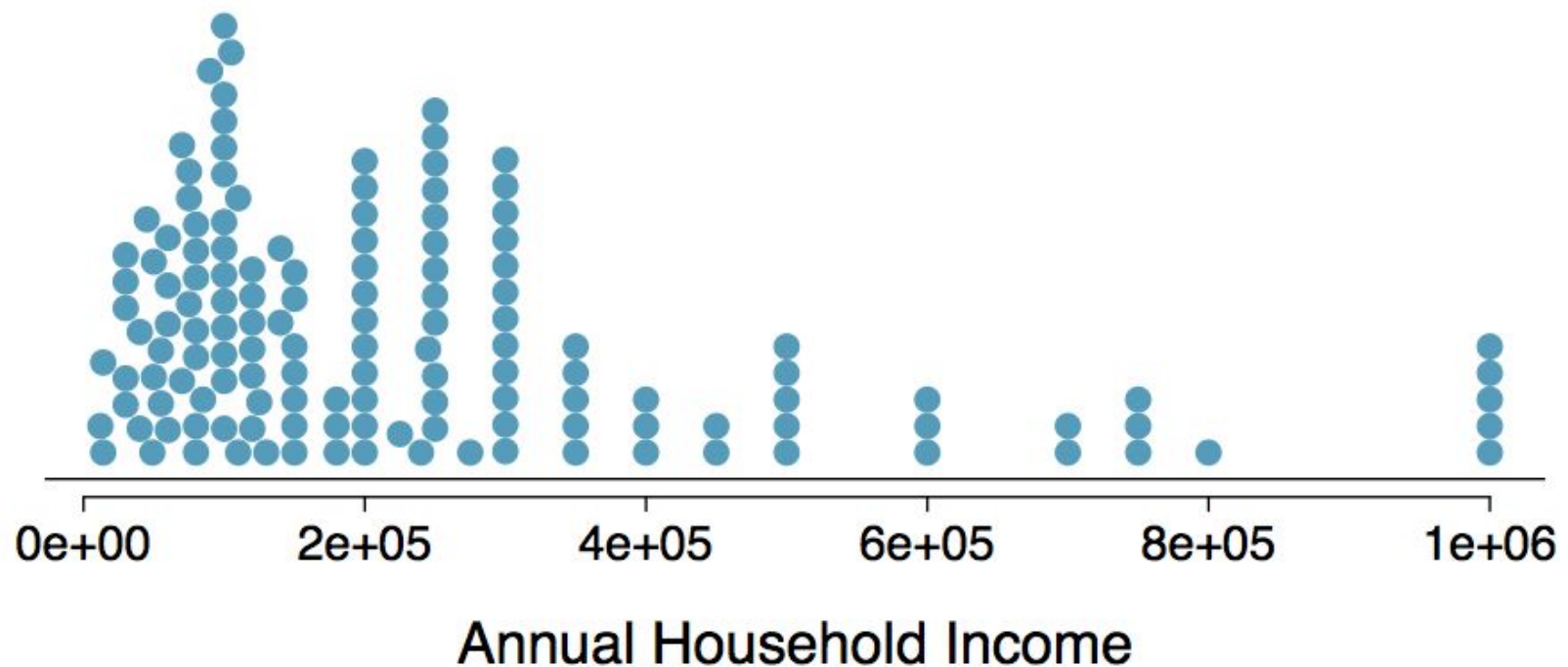
- Identify extreme skew in the distribution.
- Identify data collection and entry errors.
- Provide insight into interesting features of the data.

# Extreme Observations

How would sample statistics such as mean, median, SD, and IQR of household income be affected if the largest value was replaced with \$10 million? What if the smallest value was replaced with \$10 million?



# Robust Statistics



| scenario                      | robust |      | not robust |      |
|-------------------------------|--------|------|------------|------|
|                               | median | IQR  | $\bar{x}$  | $s$  |
| original data                 | 190K   | 200K | 245K       | 226K |
| move largest to \$10 million  | 190K   | 200K | 309K       | 853K |
| move smallest to \$10 million | 200K   | 200K | 316K       | 854K |

# Robust Statistics

Median and IQR are more robust to skewness and outliers than mean and SD. Therefore,

- for skewed distributions it is often more helpful to use median and IQR to describe the center and spread
- for symmetric distributions it is often more helpful to use the mean and SD to describe the center and spread

# Robust Statistics

Median and IQR are more robust to skewness and outliers than mean and SD. Therefore,

- for skewed distributions it is often more helpful to use median and IQR to describe the center and spread
- for symmetric distributions it is often more helpful to use the mean and SD to describe the center and spread

If you would like to estimate the typical household income for a student, would you be more interested in the mean or median income?

# Robust Statistics

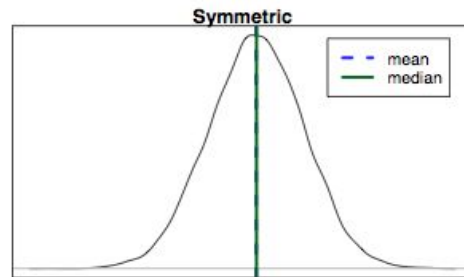
Median and IQR are more robust to skewness and outliers than mean and SD. Therefore,

- for skewed distributions it is often more helpful to use median and IQR to describe the center and spread
- for symmetric distributions it is often more helpful to use the mean and SD to describe the center and spread



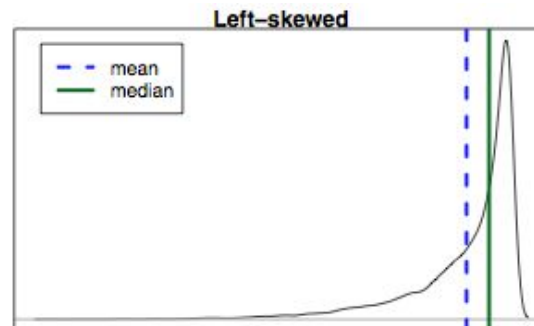
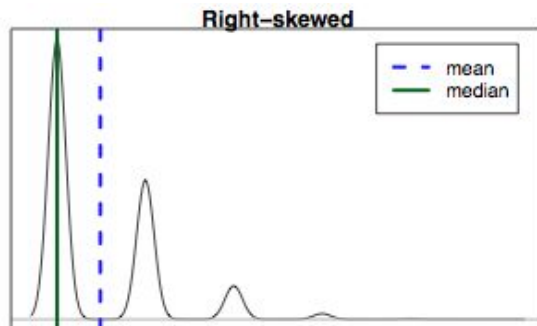
# Mean vs. Median

If the distribution is symmetric, center is often defined as the mean:  
mean  $\sim$  median



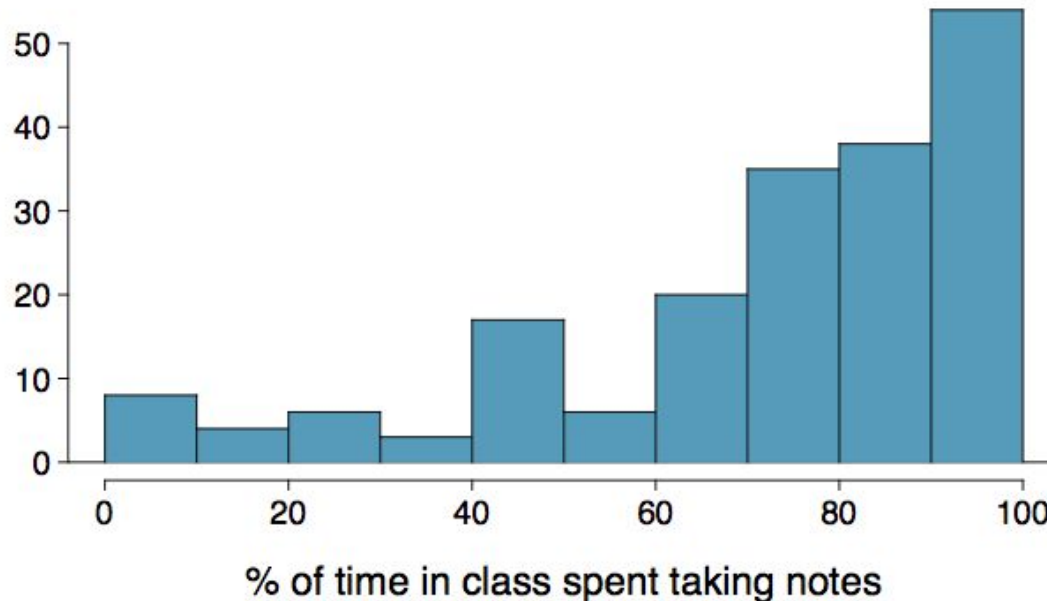
If the distribution is skewed or has extreme outliers, center is often defined as the median

- Right-skewed: mean  $>$  median
- Left-skewed: mean  $<$  median



# Practice

Which is most likely true for the distribution of percentage of time actually spent taking notes in class versus on Facebook, Twitter, etc.?



(a) mean > median

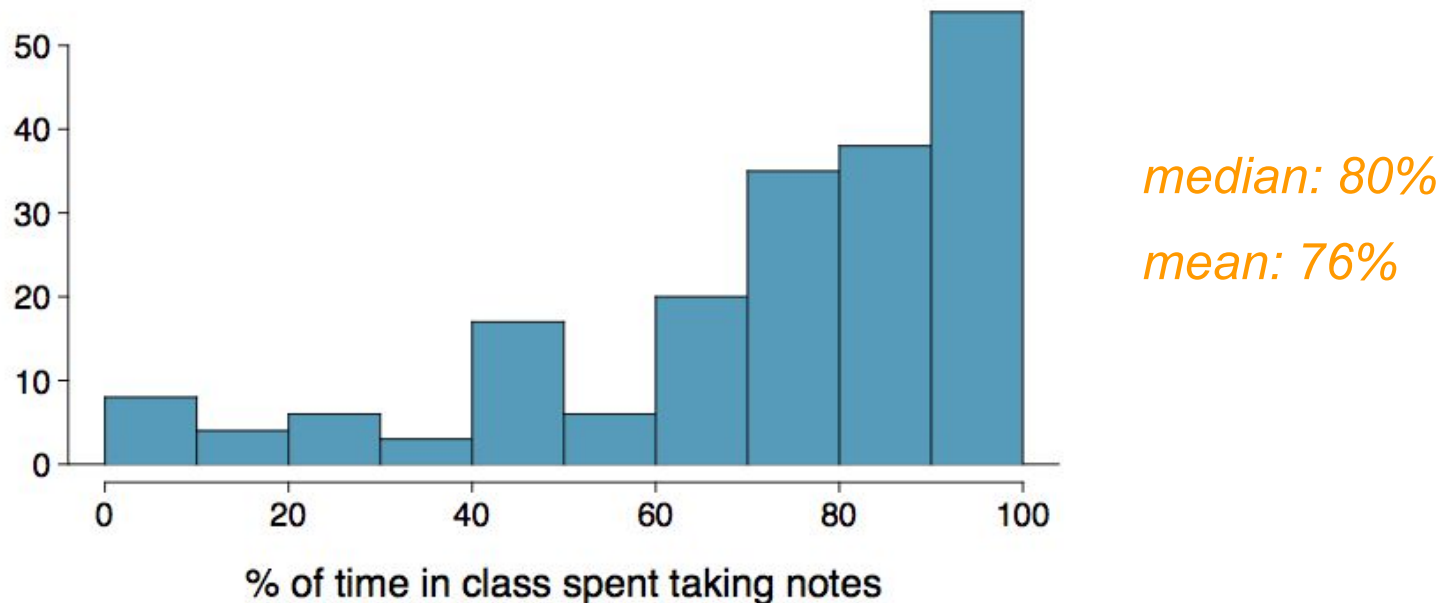
(b) mean ~ median

(c) mean < median

(d) impossible to tell

# Practice

Which is most likely true for the distribution of percentage of time actually spent taking notes in class versus on Facebook, Twitter, etc.?



(a) mean > median

(b) mean ~ median

*(c) mean < median*

(d) impossible to tell

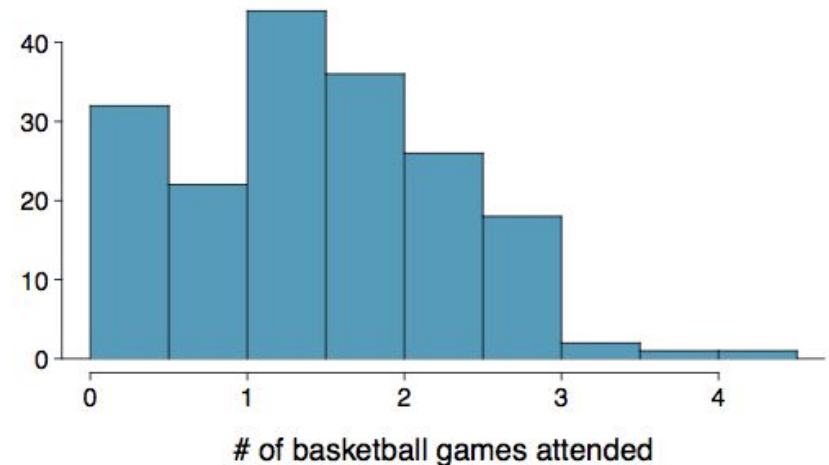
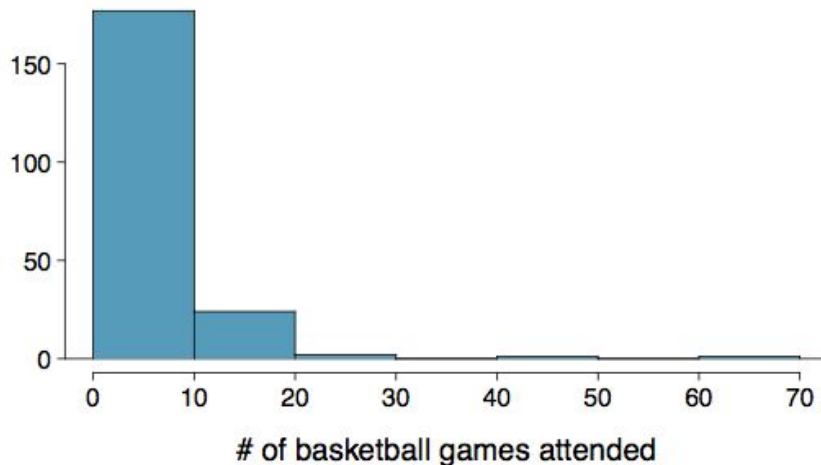
# Extremely Skewed Data

When data are extremely skewed, transforming them might make modeling easier. A common transformation is the *log transformation*.

# Extremely Skewed Data

When data are extremely skewed, transforming them might make modeling easier. A common transformation is the [log transformation](#).

The histograms on the left shows the distribution of number of basketball games attended by students. The histogram on the right shows the distribution of log of number of games attended.



# Pros and Cons of Transformations

- Skewed data are easier to model with when they are transformed because outliers tend to become far less prominent after an appropriate transformation.

|            |    |    |    |     |
|------------|----|----|----|-----|
| # of games | 70 | 50 | 25 | ... |
|------------|----|----|----|-----|

|            |      |      |      |     |
|------------|------|------|------|-----|
| # of games | 4.25 | 3.91 | 3.22 | ... |
|------------|------|------|------|-----|

- However, results of an analysis might be difficult to interpret because the log of a measured variable is usually meaningless.

# Pros and Cons of Transformations

- Skewed data are easier to model with when they are transformed because outliers tend to become far less prominent after an appropriate transformation.

|            |    |    |    |     |
|------------|----|----|----|-----|
| # of games | 70 | 50 | 25 | ... |
|------------|----|----|----|-----|

|            |      |      |      |     |
|------------|------|------|------|-----|
| # of games | 4.25 | 3.91 | 3.22 | ... |
|------------|------|------|------|-----|

- However, results of an analysis might be difficult to interpret because the log of a measured variable is usually meaningless.

What other variables would you expect to be extremely skewed?

# Pros and Cons of Transformations

- Skewed data are easier to model with when they are transformed because outliers tend to become far less prominent after an appropriate transformation.

|            |    |    |    |     |
|------------|----|----|----|-----|
| # of games | 70 | 50 | 25 | ... |
|------------|----|----|----|-----|

|            |      |      |      |     |
|------------|------|------|------|-----|
| # of games | 4.25 | 3.91 | 3.22 | ... |
|------------|------|------|------|-----|

- However, results of an analysis might be difficult to interpret because the log of a measured variable is usually meaningless.

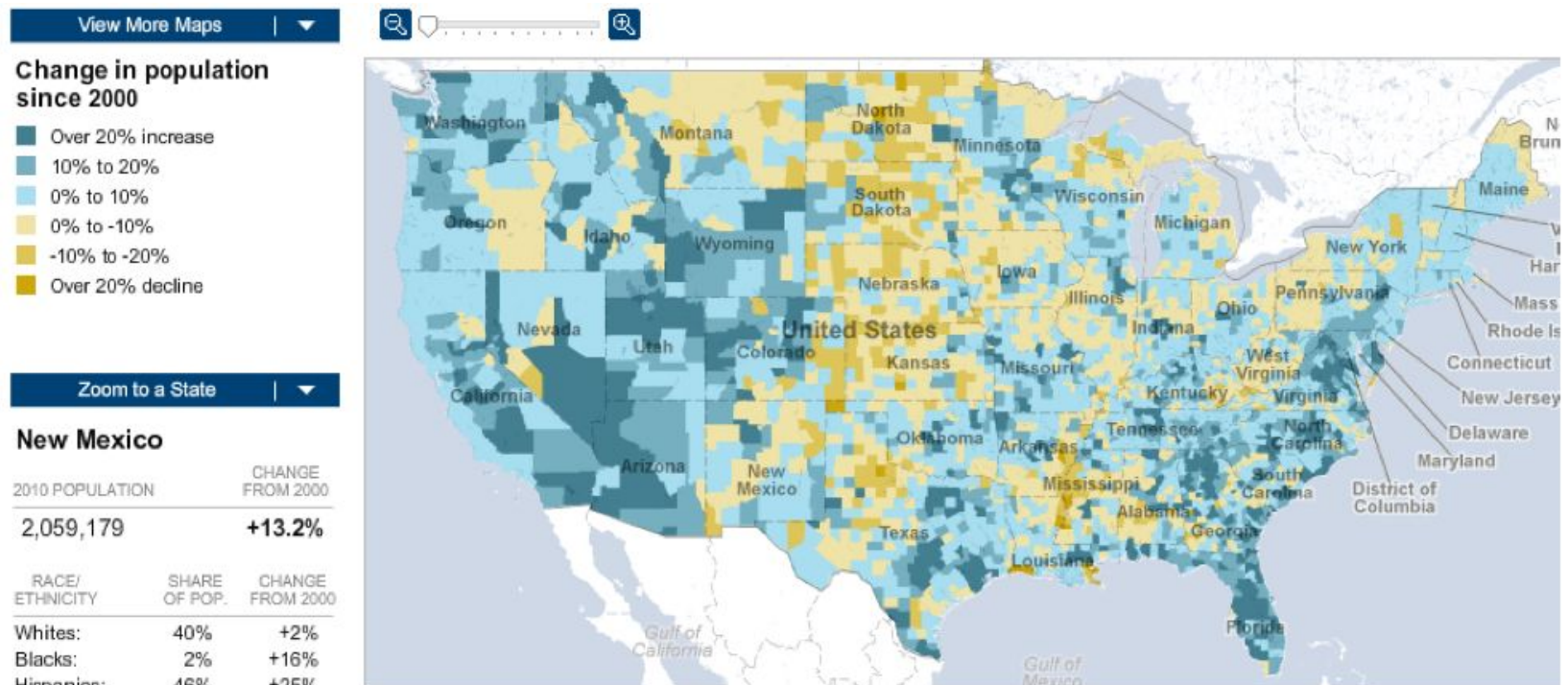
What other variables would you expect to be extremely skewed?

*Salary, housing prices, etc.*



# Intensity Maps

What patterns are apparent in the change in population between 2000 and 2010?



<http://projects.nytimes.com/census/2010/map>