

# Introduction to Statistical Inference

## Lecture 04: Data Basics

Mohammad-Reza A. Dehaqani

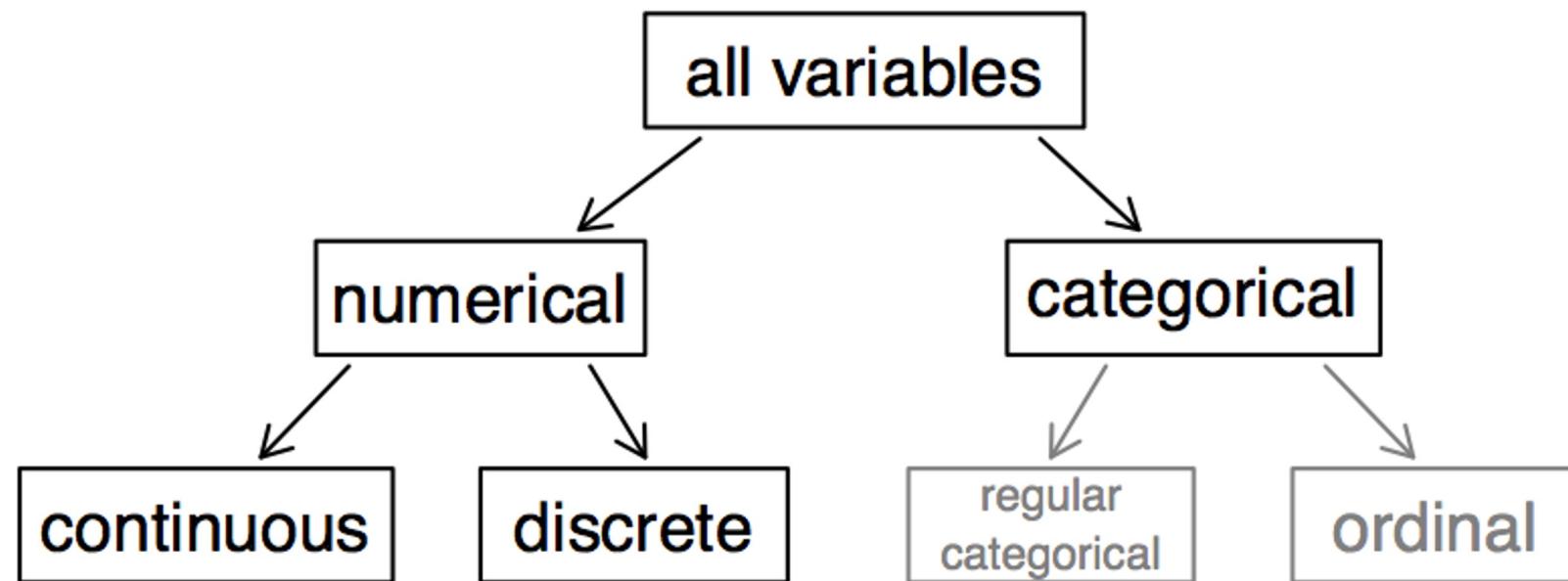
[dehaqani@ut.ac.ir](mailto:dehaqani@ut.ac.ir)

# Data matrix

Data collected on students in a statistics class on a variety of variables:

		<i>variable</i>			
			↓		
Stu.	gender	intro_extra	...	dread	
1	male	extravert	...	3	
2	female	extravert	...	2	
3	female	introvert	...	4	←
4	female	extravert	...	2	<i>observation</i>
:	:	:	:	:	
86	male	extravert	...	3	

# Types of variables



# Types of variables (cont.)

	gender	sleep	bedtime	countries	dread
1	male	5	12-2	13	3
2	female	7	10-12	7	2
3	female	5.5	12-2	1	4
4	female	7	12-2		2
5	female	3	12-2	1	3
6	female	3	12-2	9	4

- gender: *categorical*
- sleep: *numerical, continuous*
- bedtime: *categorical, ordinal*
- countries: *numerical, discrete*
- dread: *categorical, ordinal - could also be used as numerical*

# Populations and Samples

PHYS ED | AUGUST 29, 2012, 12:01 AM | 21 Comments

## Finding Your Ideal Running Form

By GRETCHEN REYNOLDS



David De Lossy/Getty Images

<http://well.blogs.nytimes.com/2012/08/29/finding-your-ideal-running-form>

*Sample:* Group of adult women who recently joined a running group

*Population to which results can be generalized:* Adult women, if the data are randomly sampled

# Census

- Wouldn't it be better to just include everyone and "sample" the entire population?
  - This is called a *census*.
- There are problems with taking a census:
  - It can be difficult to complete a census: there always seem to be some individuals who are hard to locate or hard to measure. *And these difficult-to-find people may have certain characteristics that distinguish them from the rest of the population.*
  - Populations rarely stand still. Even if you could take a census, the population changes constantly, so it's never possible to get a perfect measure.
  - Taking a census may be more complex than sampling.

# Exploratory analysis to inference

- Sampling is natural.
- Think about sampling something you are cooking - you taste (examine) a small part of what you're cooking to get an idea about the dish as a whole.
- When you taste a spoonful of soup and decide the spoonful you tasted isn't salty enough, that's *exploratory analysis*.
- If you generalize and conclude that your entire soup needs salt, that's an *inference*.
- For your inference to be valid, the spoonful you tasted (the sample) needs to be *representative* of the entire pot (the population).
  - If your spoonful comes only from the surface and the salt is collected at the bottom of the pot, what you tasted is probably not representative of the whole pot.
  - If you first stir the soup thoroughly before you taste, your spoonful will more likely be representative of the whole pot.

# Sampling bias

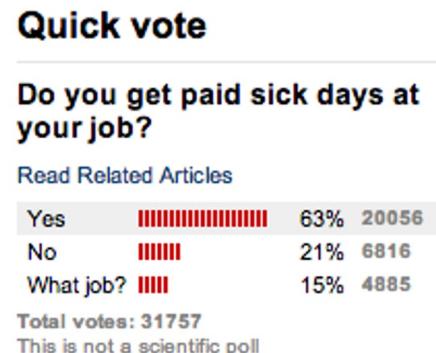
- **Non-response:** If only a small fraction of the randomly sampled people choose to respond to a survey, the sample may no longer be representative of the population.
- **Voluntary response:** Occurs when the sample consists of people who volunteer to respond because they have strong opinions on the issue. Such a sample will also not be representative of the population.

**Quick vote**

**Do you get paid sick days at your job?**

Yes       No  
 What job?

**VOTE** or view results



- **Convenience sample:** Individuals who are easily accessible are more likely to be included in the sample.

# Observational studies

- Researchers collect data in a way that does not directly interfere with how the data arise.
- Results of an observational study can generally be used to establish an association between the explanatory and response variables.

# Prospective vs. Retrospective Studies

A **prospective study** identifies individuals and collects information as events unfold.

- Example: The Nurses Health Study has been recruiting registered nurses and then collecting data from them using questionnaires since 1976.

**Retrospective studies** collect data after events have taken place.

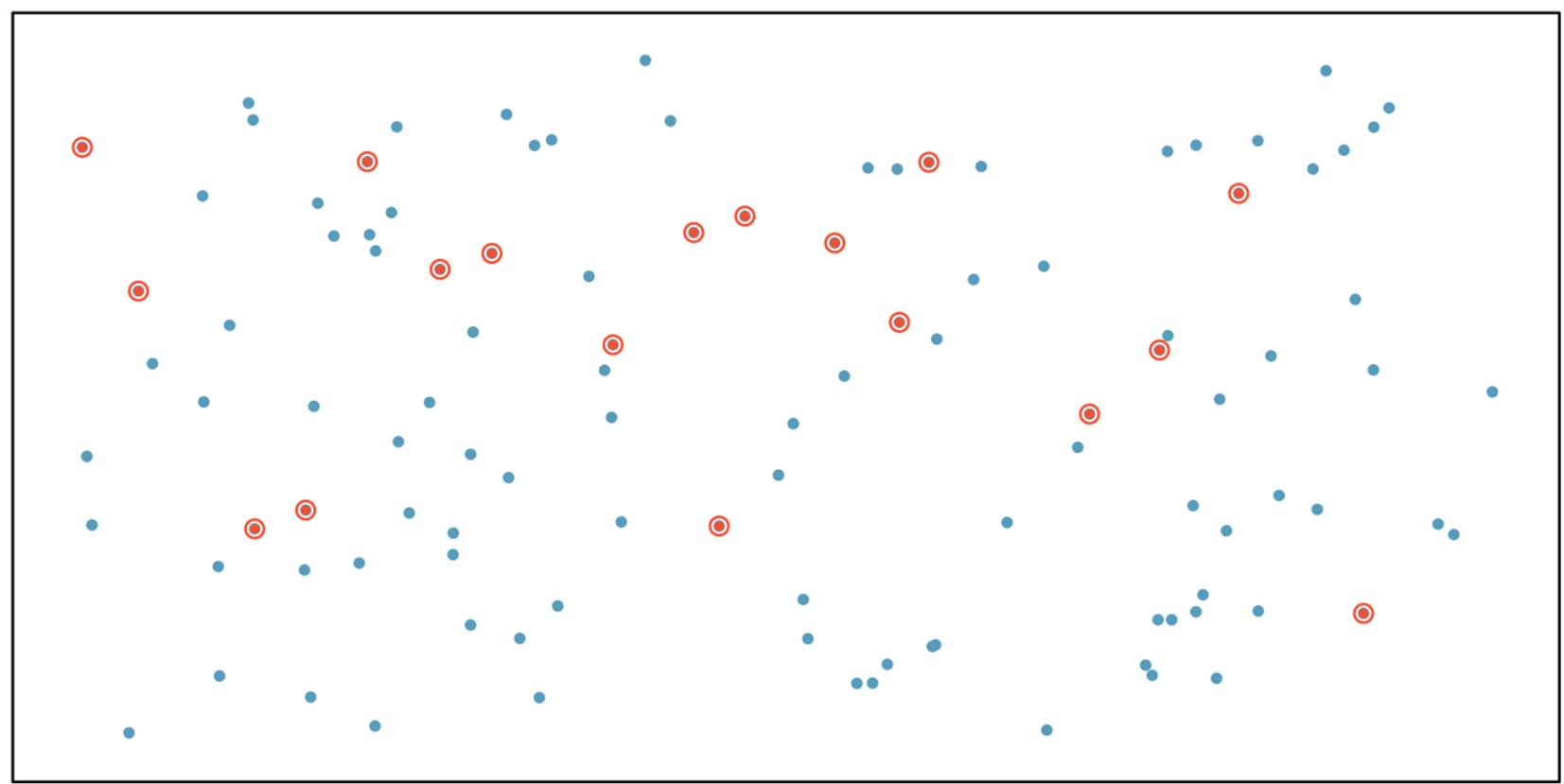
- Example: Researchers reviewing past events in medical records.

# Obtaining Good Samples

- Almost all statistical methods are based on the notion of implied randomness.
  - If observational data are not collected in a random framework from a population, these statistical methods -- the estimates and errors associated with the estimates -- are not reliable.
  - Most commonly used random sampling techniques are *simple*, *stratified*, and *cluster* sampling.

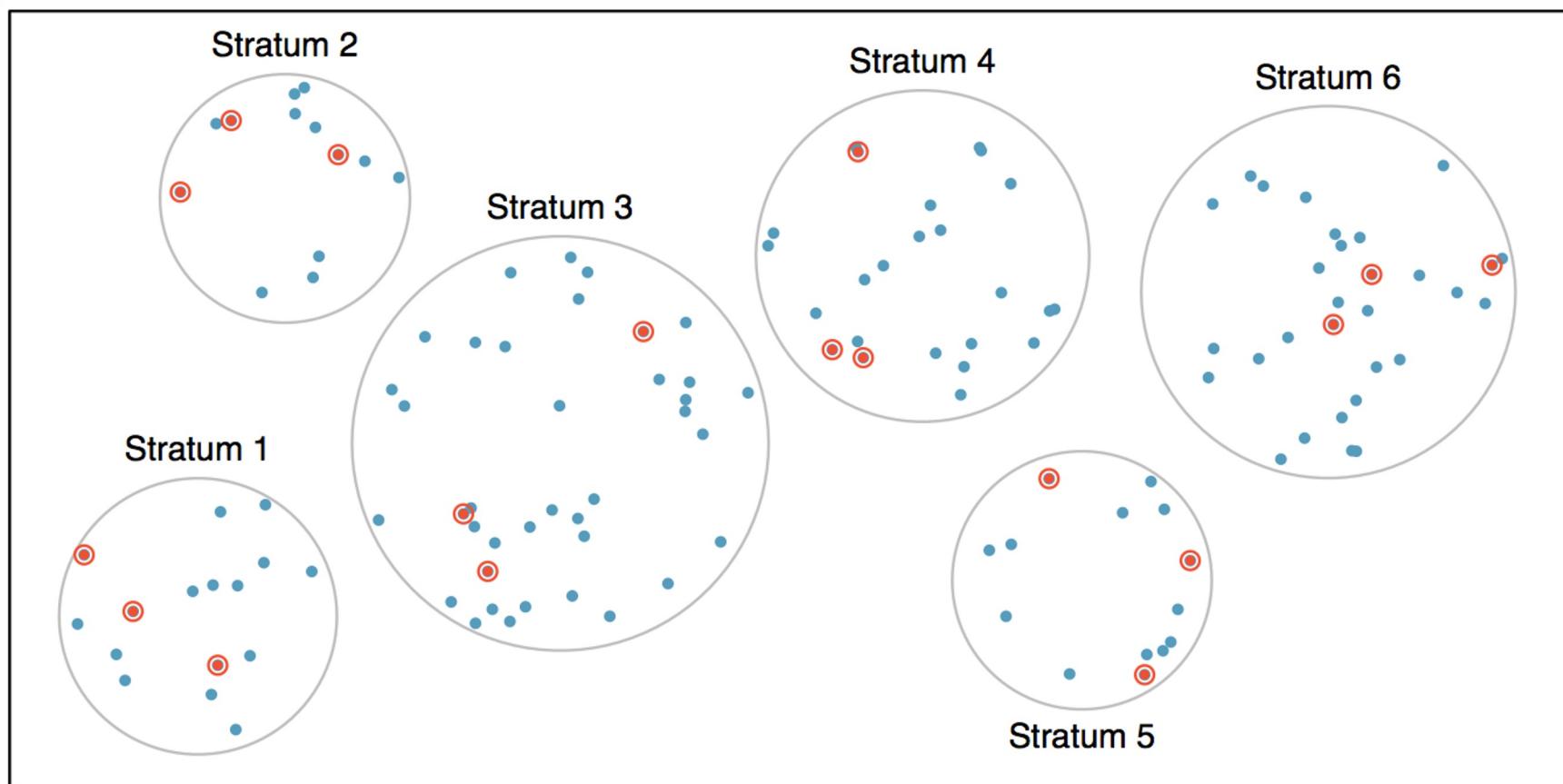
# Simple Random Sample

Randomly select cases from the population, where there is no implied connection between the points that are selected.



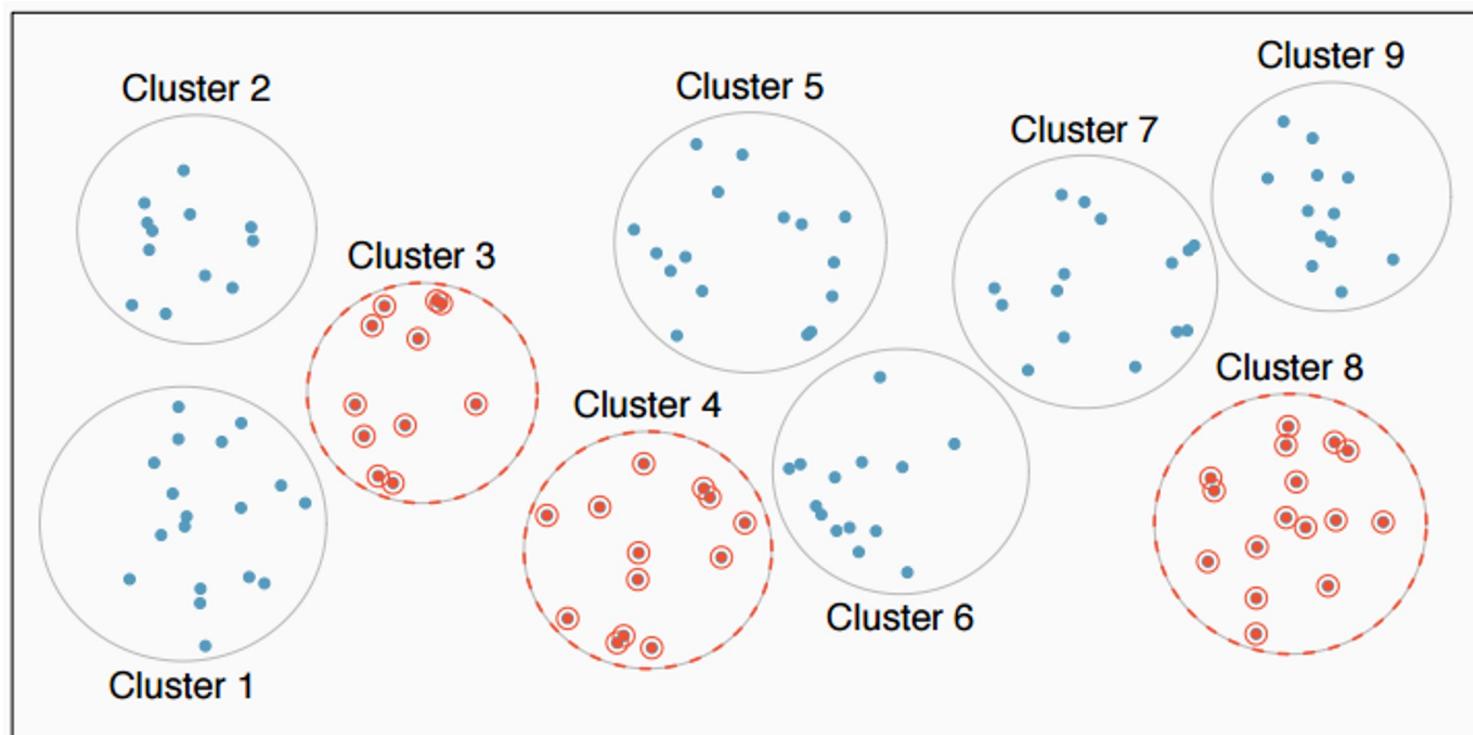
# Stratified Sample

*Strata* are made up of similar observations. We take a simple random sample from each stratum.



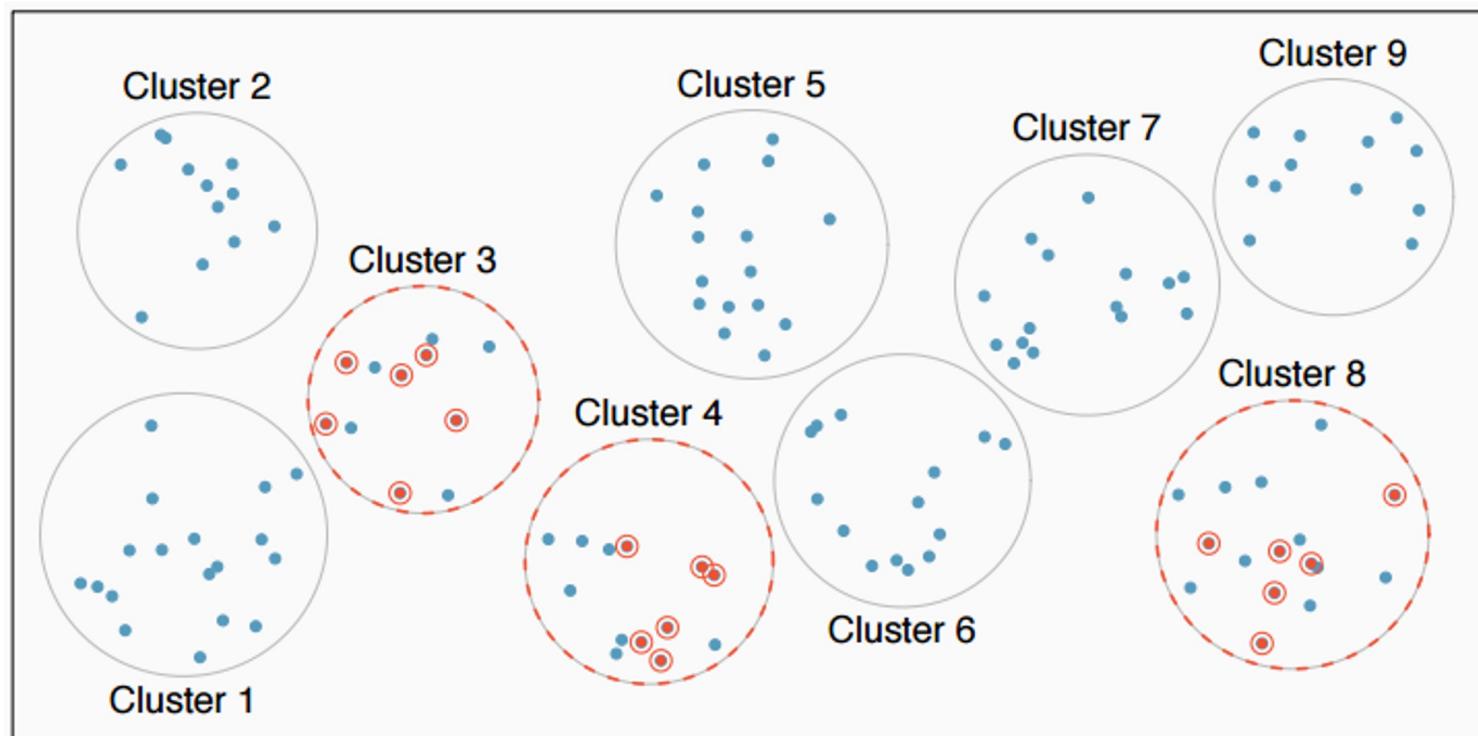
# Cluster Sample

*Clusters* are usually not made up of homogeneous observations. We take a simple random sample of clusters, and then sample all observations in that cluster. Usually preferred for economical reasons.



# Multistage Sample

*Clusters* are usually not made up of homogeneous observations. We take a simple random sample of clusters, and then take a simple random sample of observations from the sampled clusters



# Considering Numerical Data

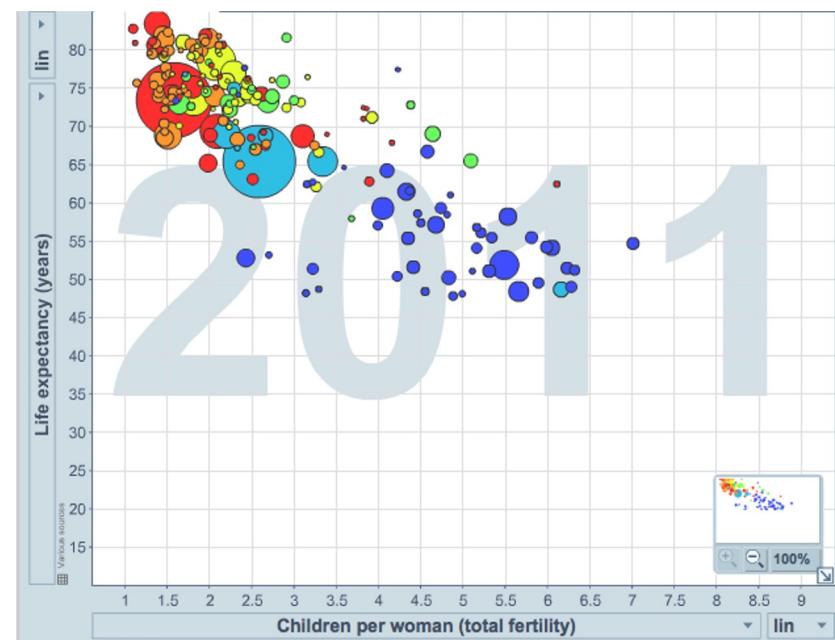
# Scatterplot

*Scatterplots* are useful for visualizing the relationship between two numerical variables.

Do life expectancy and total fertility appear to be *associated* or *independent*?

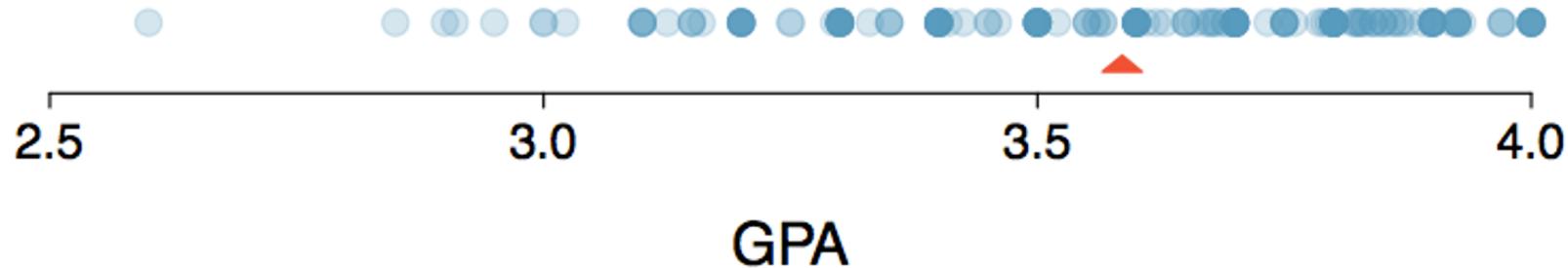
They appear to be linearly and negatively associated: as fertility increases, life expectancy decreases.

Was the relationship the same throughout the years, or did it change?  
The relationship changed over the years.



<http://www.gapminder.org/world>

# Dot Plots & Mean

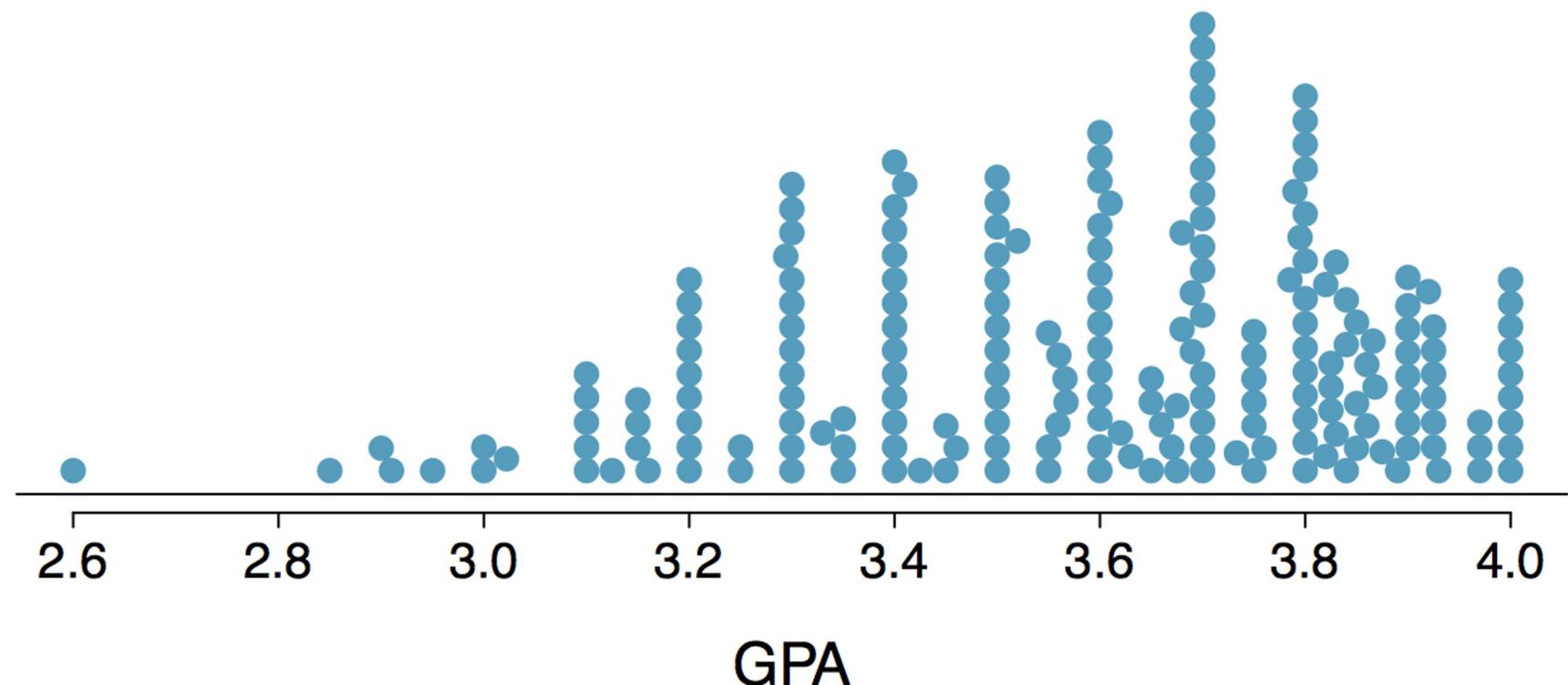


The *mean*, also called the *average* (marked with a triangle in the above plot), is one way to measure the center of a *distribution* of data.

The mean GPA is **3.59**.

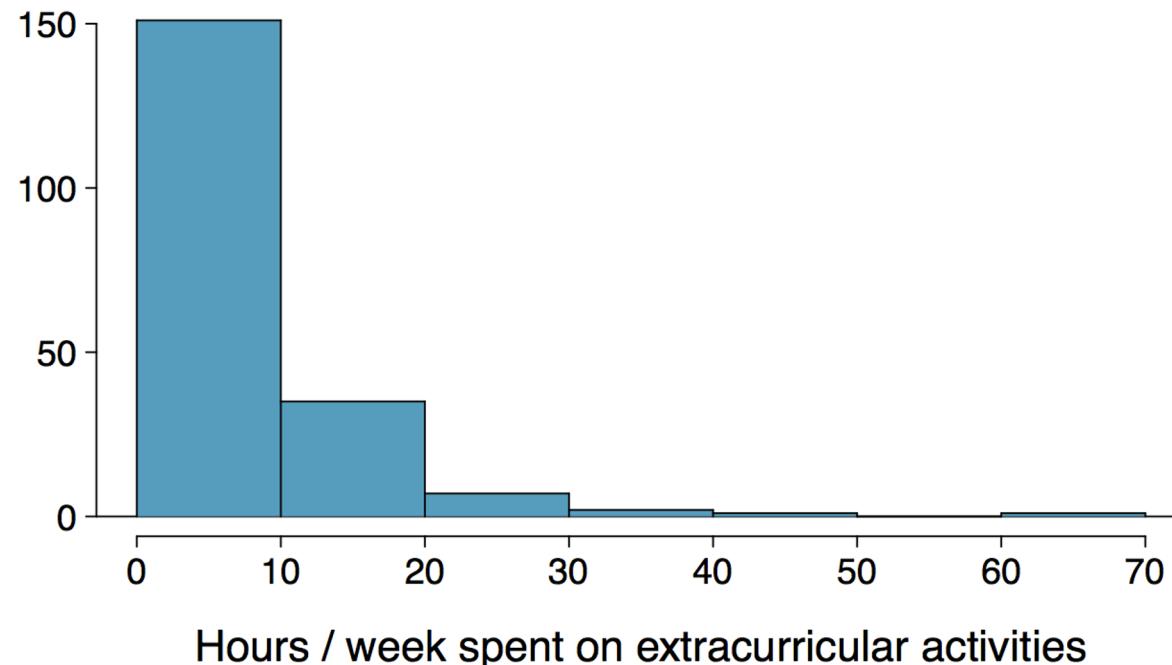
# Stacked Dot Plot

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



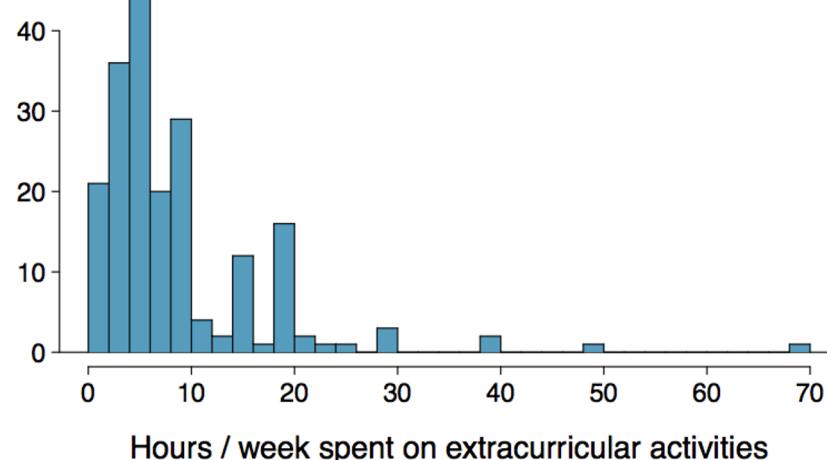
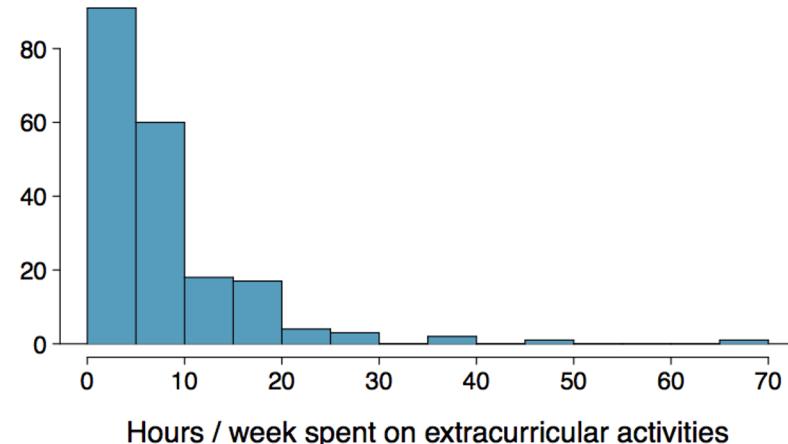
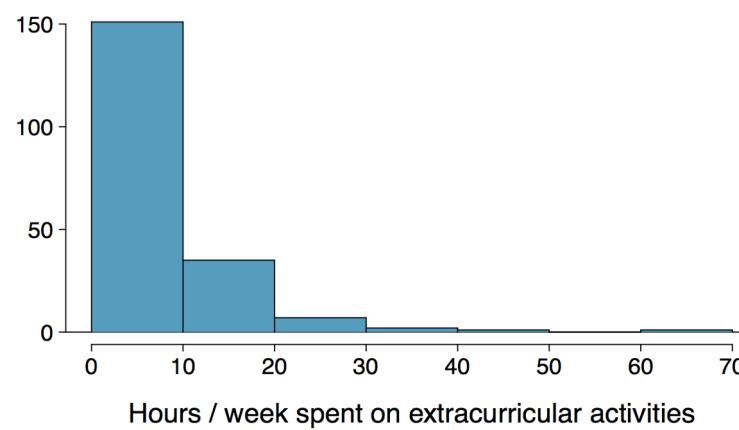
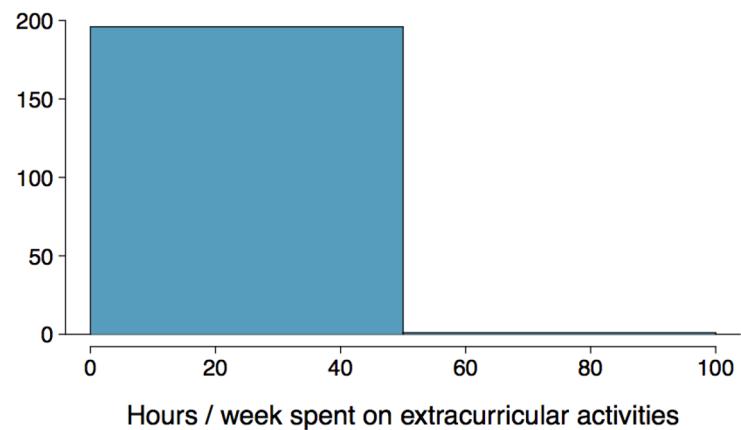
# Histograms - Extracurricular Hours

- Histograms provide a view of the *data density*. Higher bars represent where the data are relatively more common.
- Histograms are especially convenient for describing the *shape* of the data distribution.
- The chosen *bin width* can alter the story the histogram is telling.



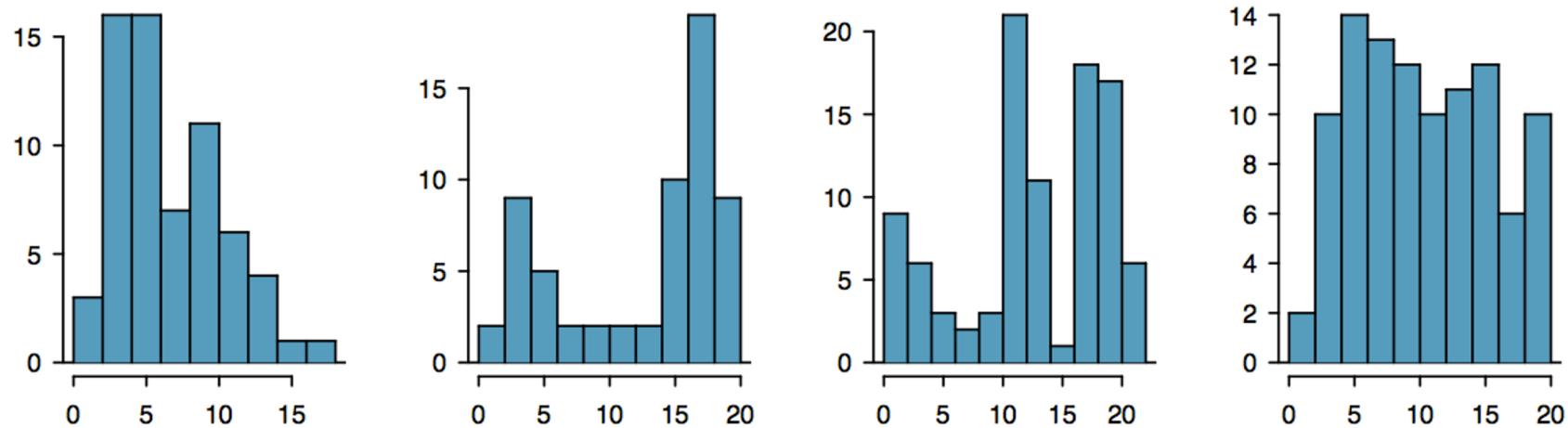
# Bin Width

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



# Shape of a Distribution: Modality

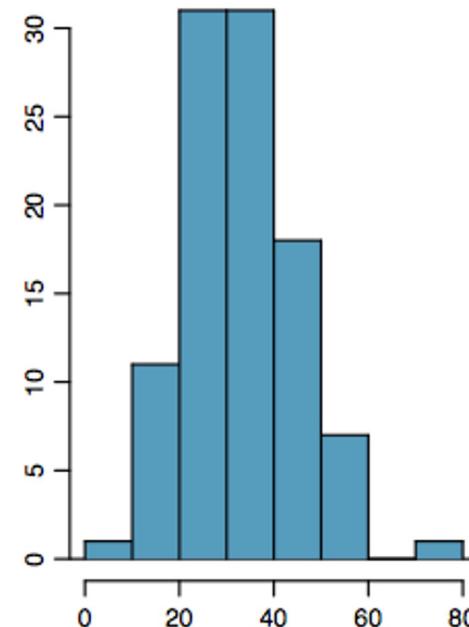
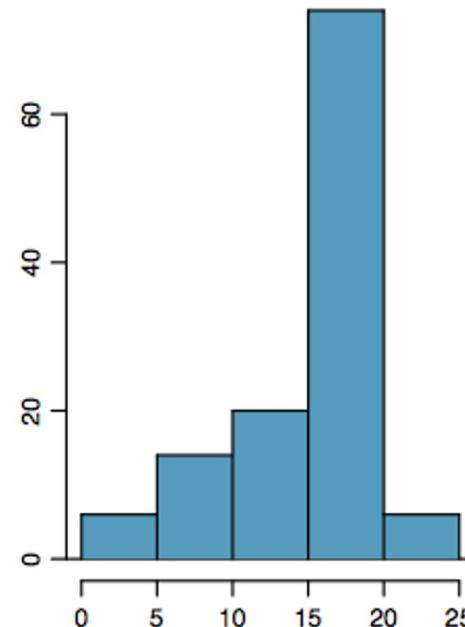
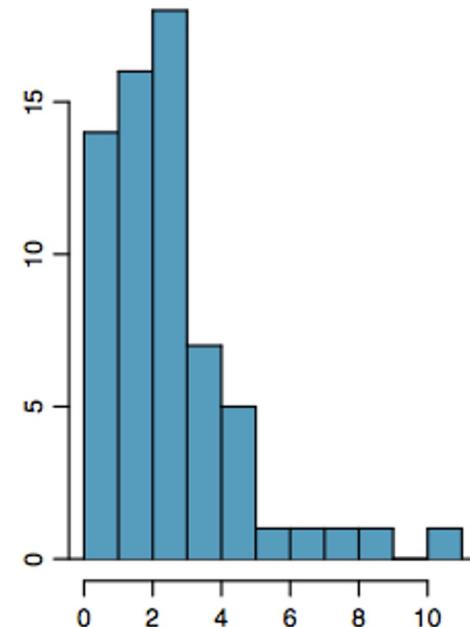
Does the histogram have a single prominent peak (*unimodal*), several prominent peaks (*bimodal/multimodal*), or no apparent peaks (*uniform*)?



**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: Skewness

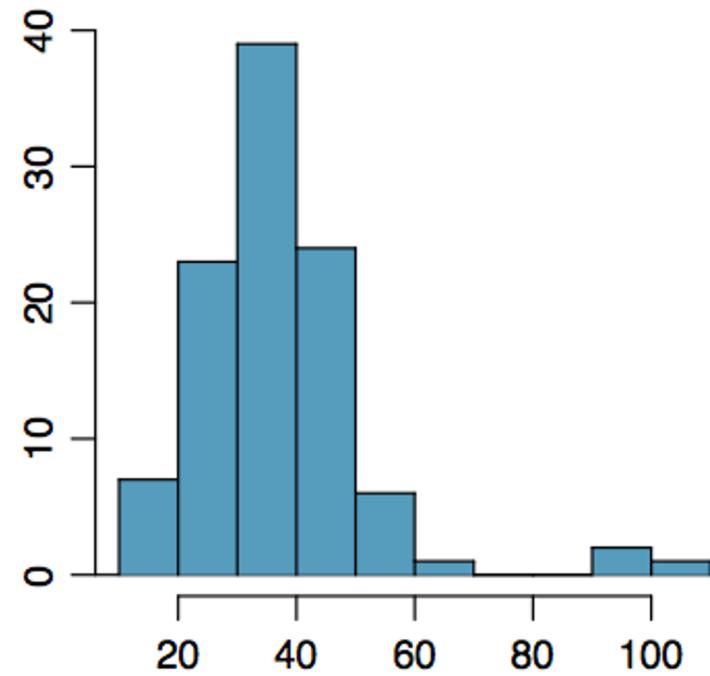
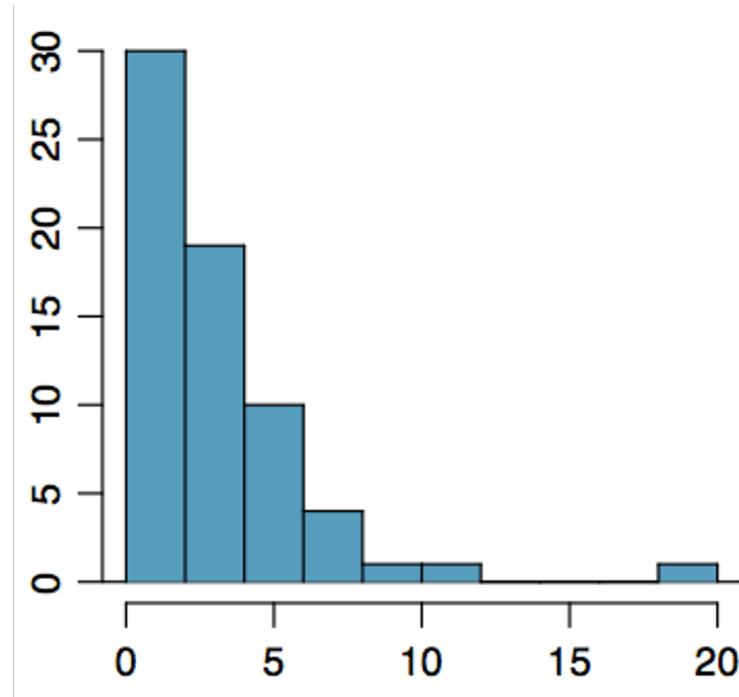
Is the histogram *right skewed*, *left skewed*, or *symmetric*?



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

# Shape of a Distribution: Unusual Observations

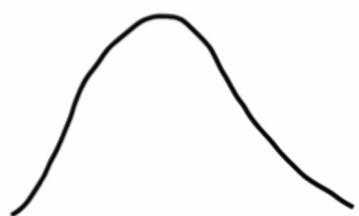
Are there any unusual observations or potential *outliers*?



# Commonly observed shapes of distributions

Modality

unimodal



bimodal



multimodal



uniform



Skewness

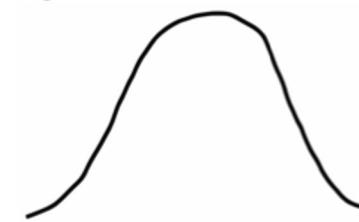
right skew



left skew



symmetric

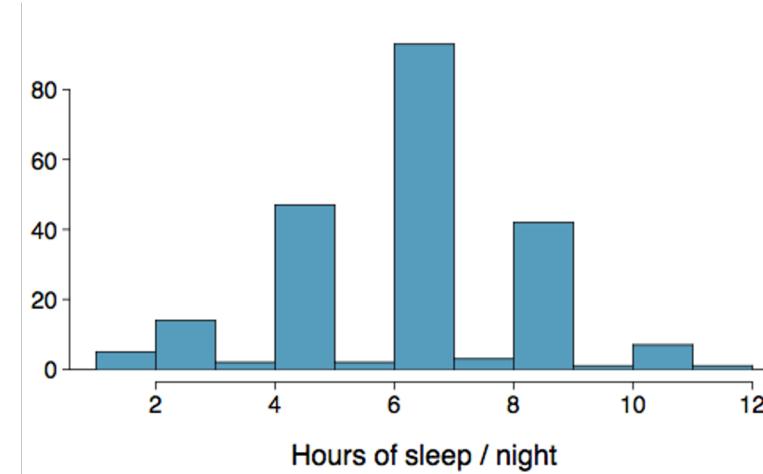


# Variance

Variance is roughly the average squared deviation from the mean.

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

- 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?

- To get rid of negatives so that observations equally distant from the mean are weighed equally.
- To weigh larger deviations more heavily.

# Standard Deviation

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

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

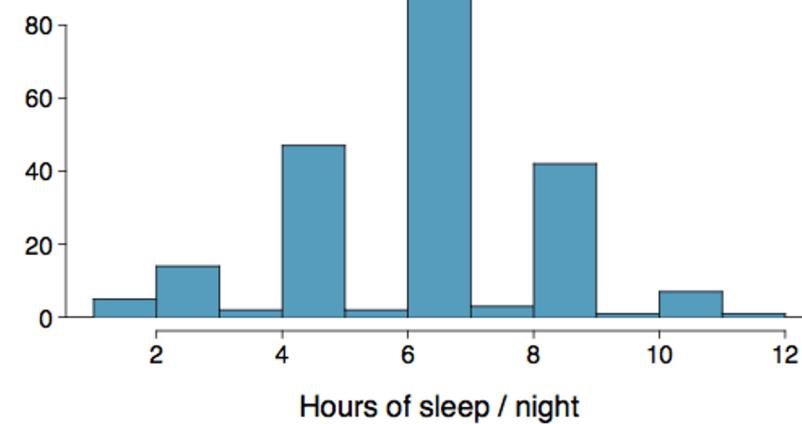
# Standard Deviation

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:

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



- 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.

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

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

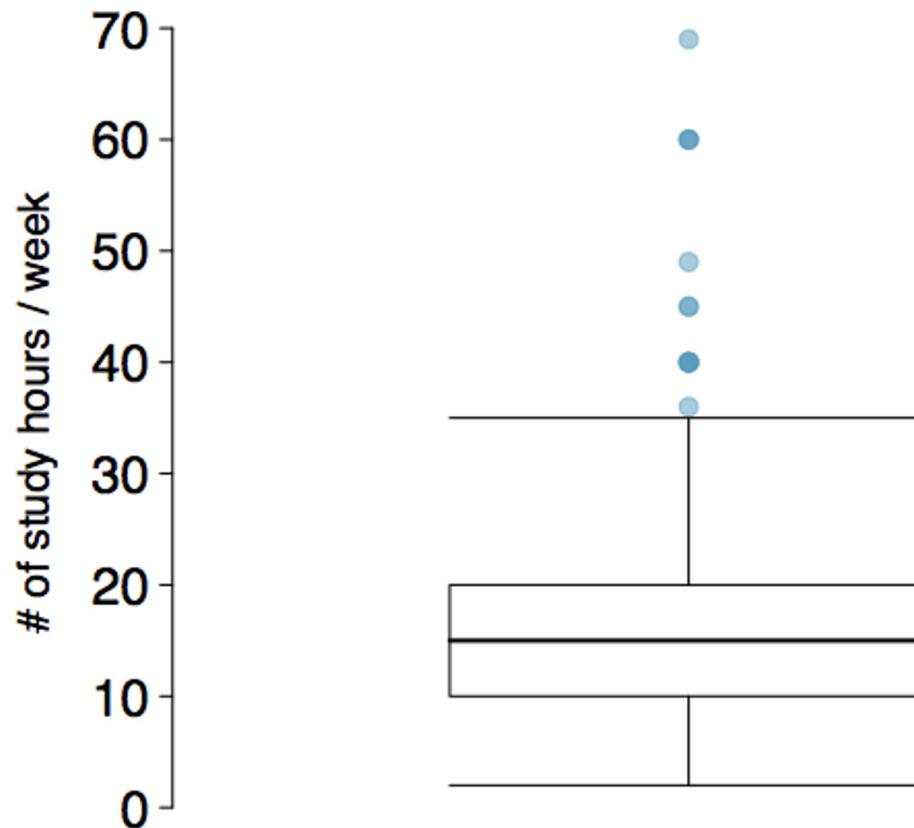
# Q1, Q3, and IQR

- The 25th percentile is also called the first quartile,  $Q_1$ .
- The 50th percentile is also called the median.
- The 75th percentile is also called the third quartile,  $Q_3$ .
- Between  $Q_1$  and  $Q_3$  is the middle 50% of the data. The range these data span is called the *interquartile range*, or the  $IQR$ .

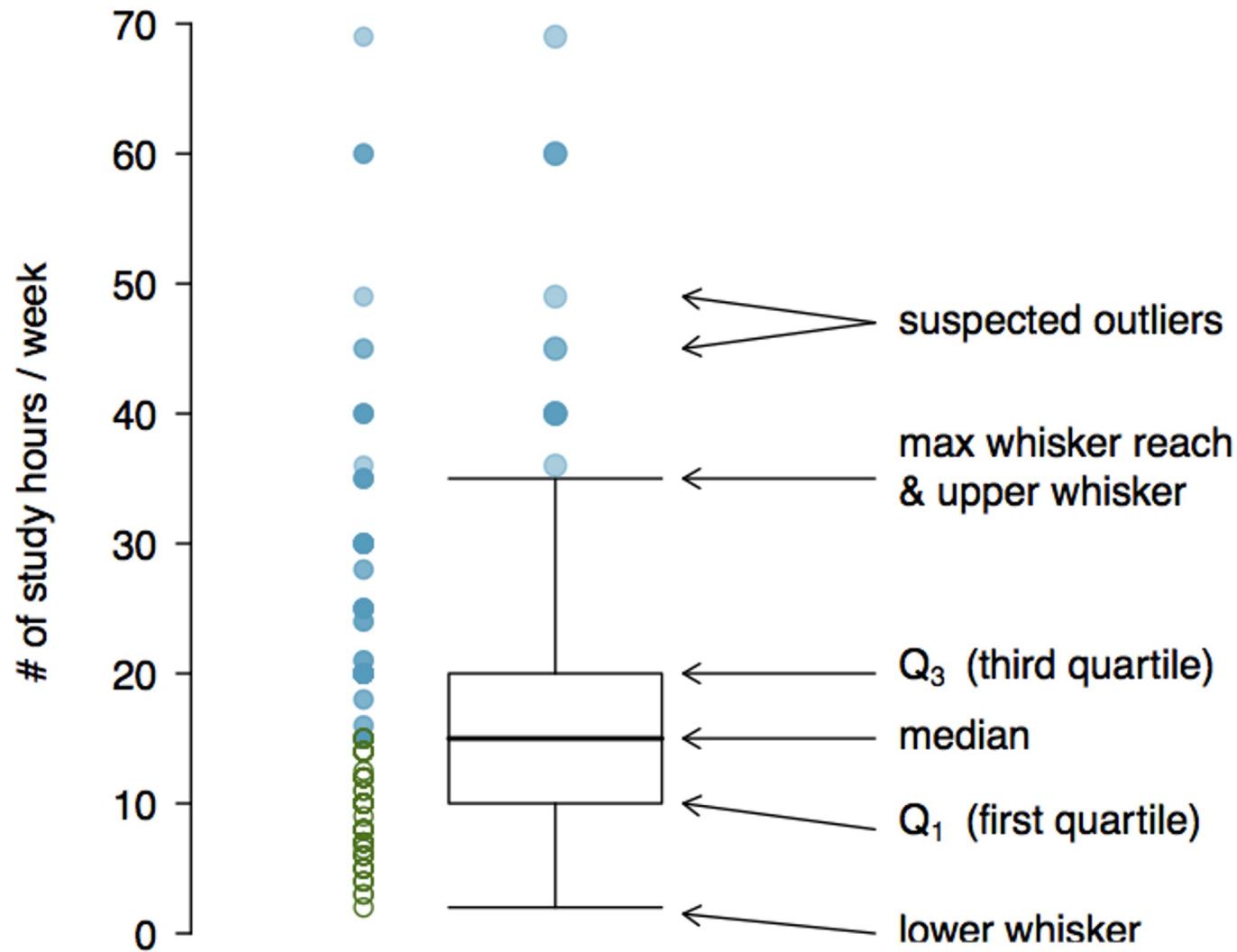
$$IQR = Q_3 - Q_1$$

# 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.



# Anatomy of a Box Plot



# 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.

# 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.

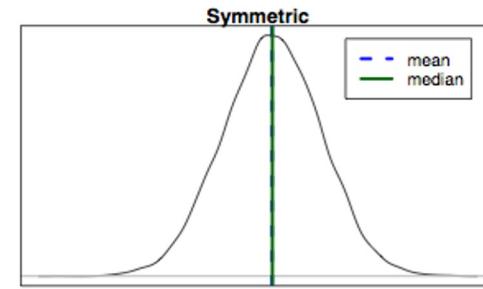
# Robust Statistics

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

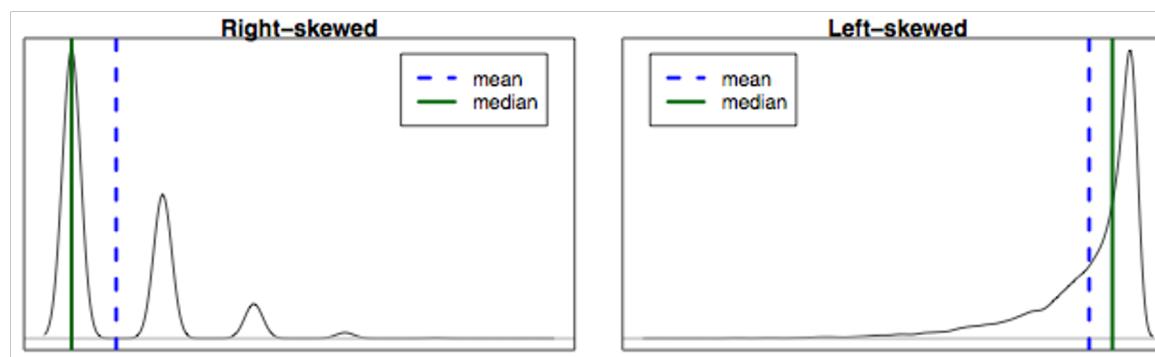
# Mean vs. Median

If the distribution is symmetric, center is often defined as the mean:  
mean ~ median



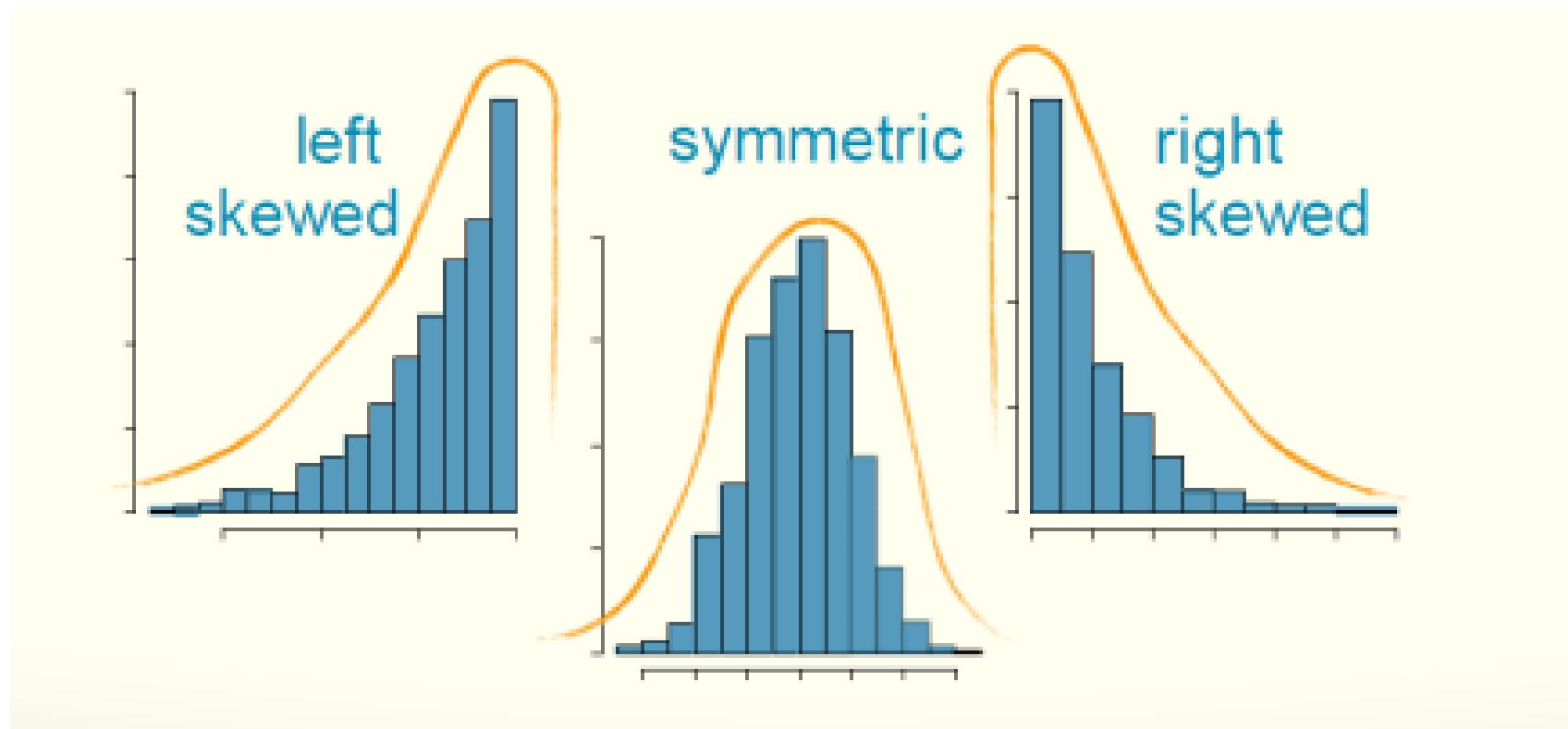
If the distribution is skewed or has extreme outliers, center is often defined as the median

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

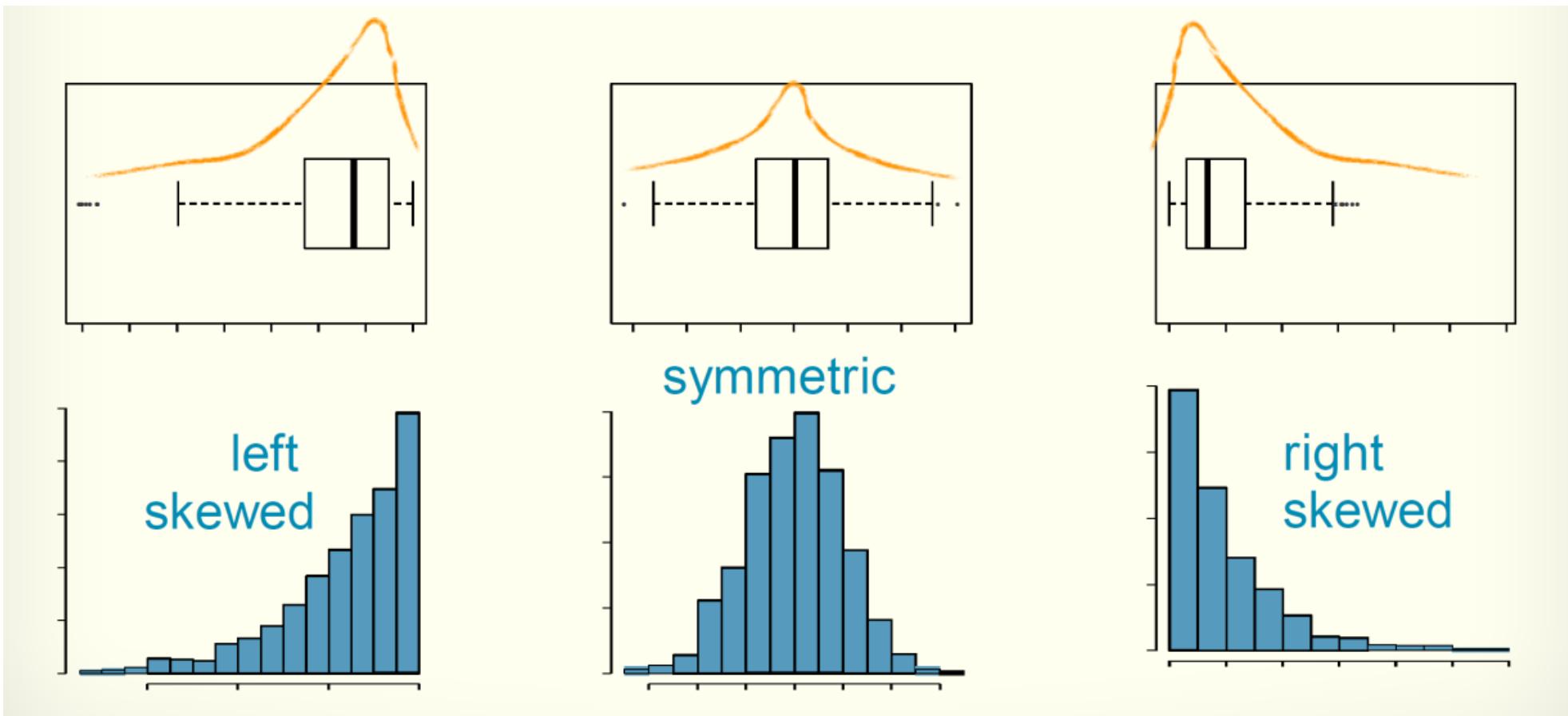


# Skewness

Distributions are skewed to the side of the long tail



# Determining the skewness from a box plot

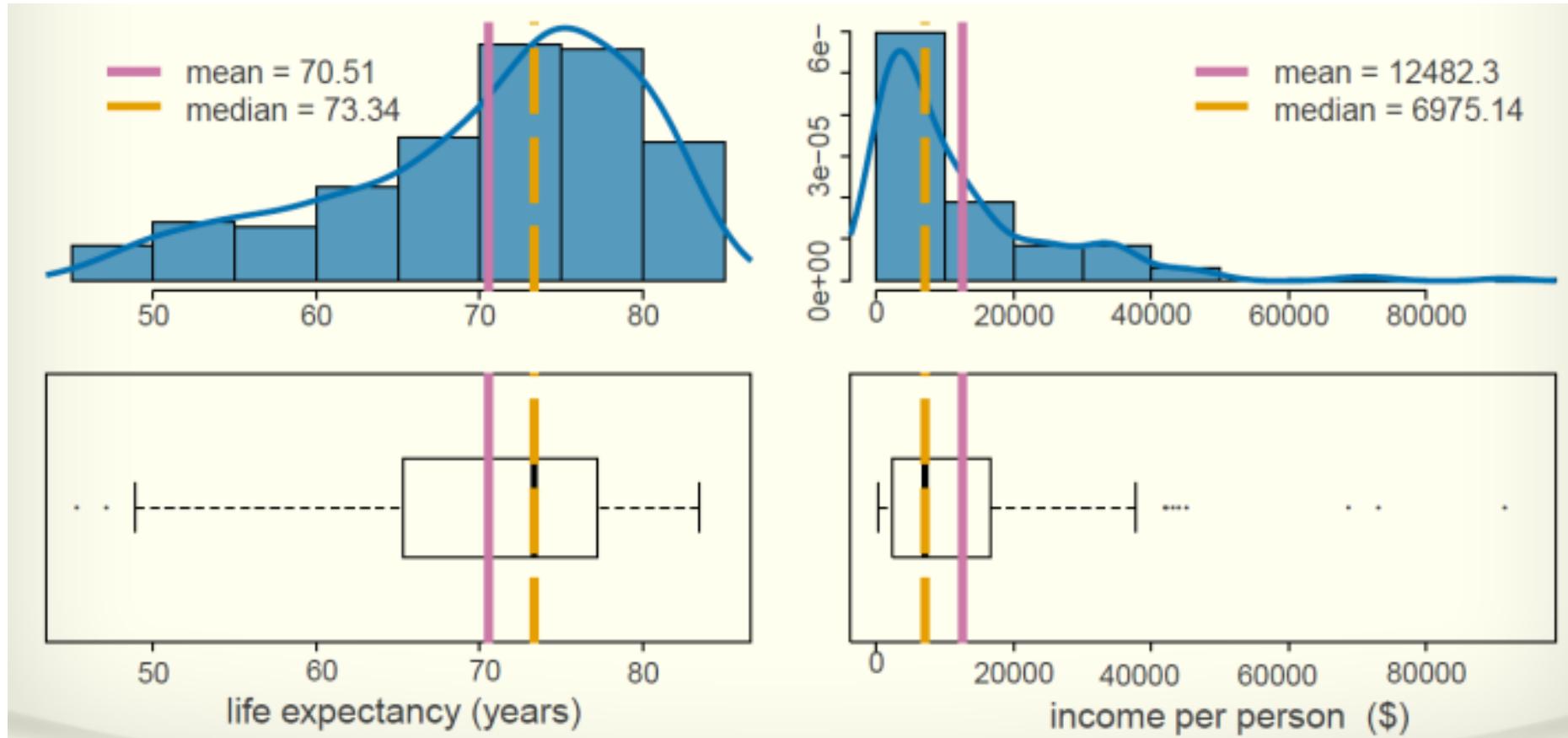


# Extremely Skewed Data

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

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

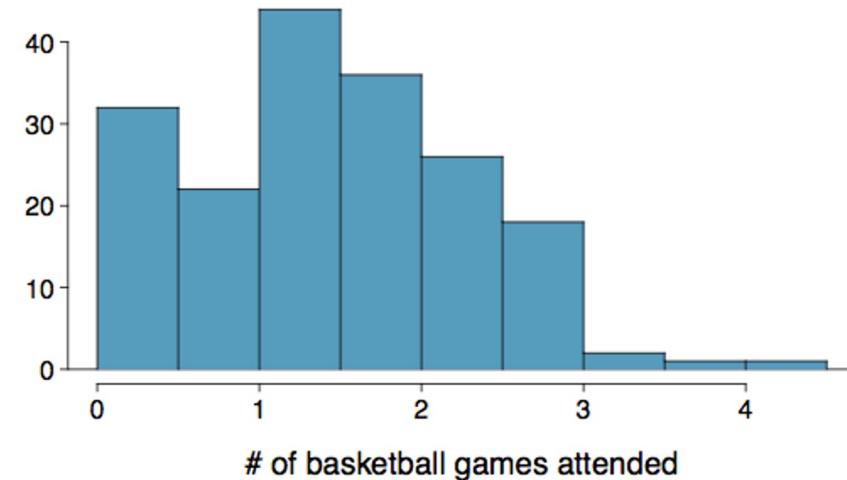
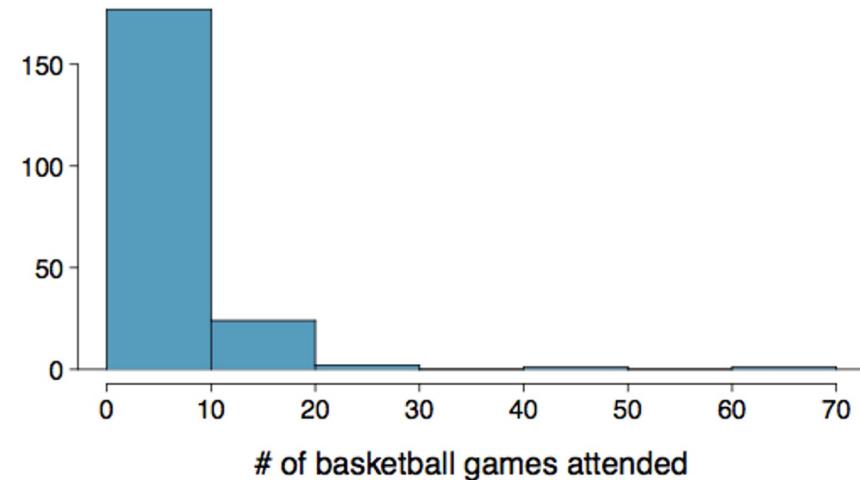
# Relation between Mean and Median



# 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.



➤ Non-parametric skewness:

$$sk = \frac{\text{mean} - \text{median}}{\text{standard deviation}} = \frac{\mu - m}{\sigma}$$

➤  $sk > 0$  : right-skewed

➤  $sk = 0$  : symmetric

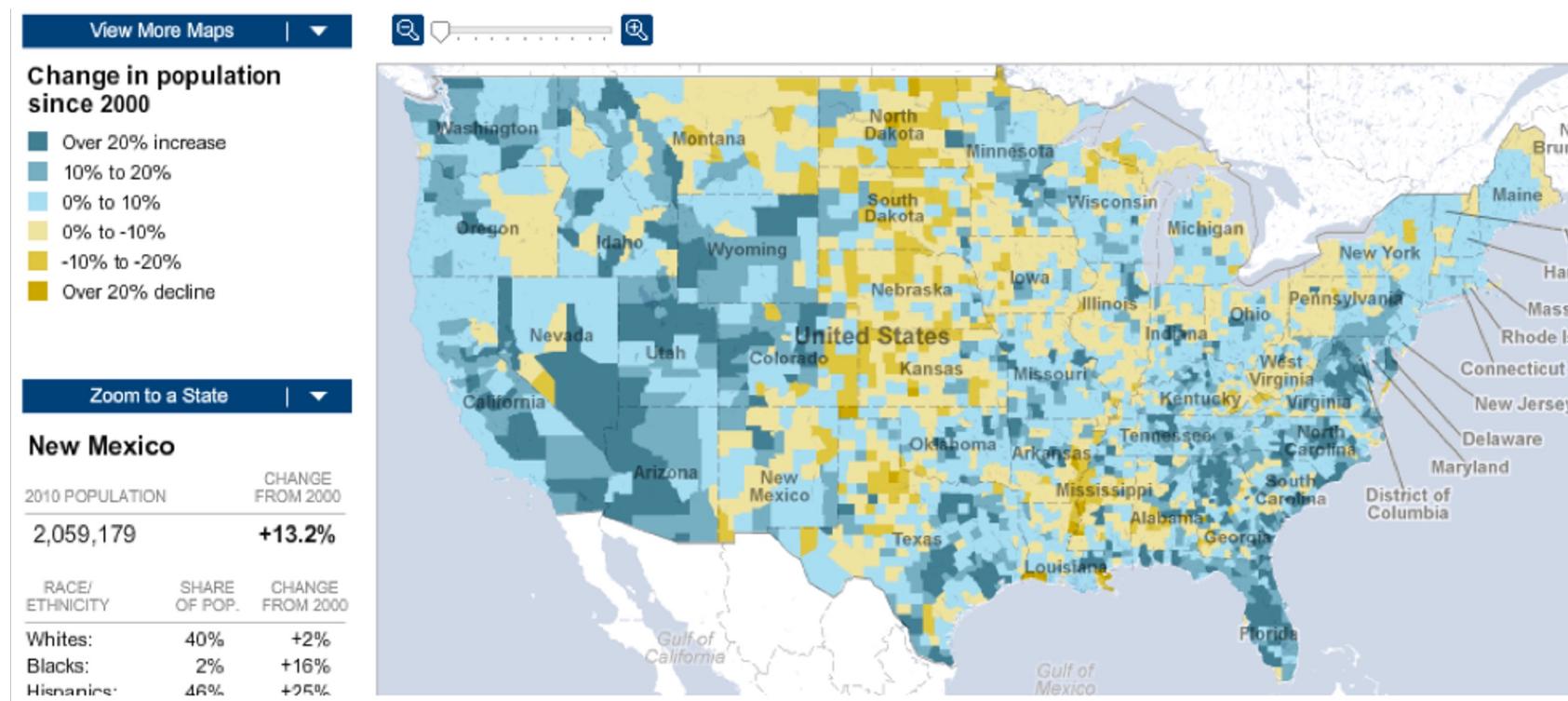
➤  $sk < 0$  : left-skewed

➤ Pearson's moment coefficient of skewness:

$$\gamma_1 = E\left(\left(\frac{X - \mu}{\sigma}\right)^3\right) = \frac{\mu_3}{\sigma^3}$$

# Intensity Maps

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



*http://projects.nytimes.com/census/2010/map*

# Considering Categorical Data

# Contingency Tables

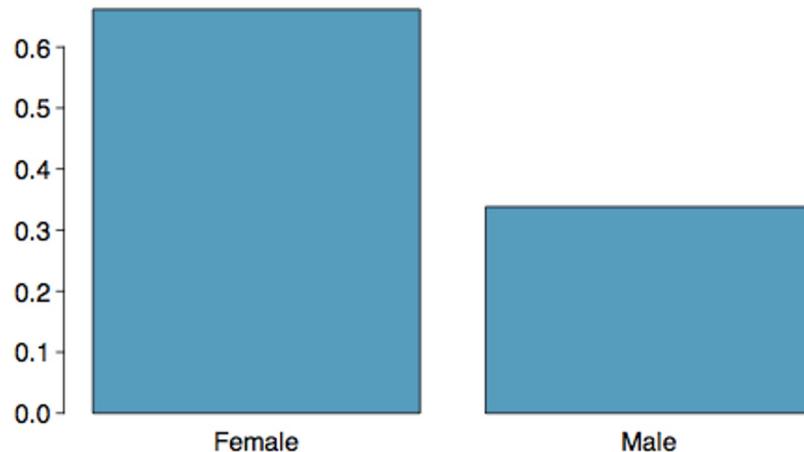
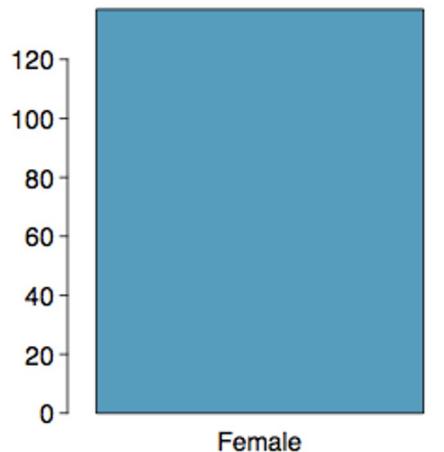
A table that summarizes data for two categorical variables is called a *contingency table*.

The contingency table below shows the distribution of students' genders and whether or not they are looking for a spouse while in college.

gender	looking for spouse		
	No	Yes	Total
Female	86	51	137
Male	52	18	70
Total	138	69	207

# Bar Plots

A *bar plot* is a common way to display a single categorical variable. A bar plot where proportions instead of frequencies are shown is called a *relative frequency bar plot*.



How are bar plots different than histograms?

Bar plots are used for displaying distributions of categorical variables, while histograms are used for numerical variables. The x-axis in a histogram is a number line, hence the order of the bars cannot be changed, while in a bar plot the categories can be listed in any order (though some orderings make more sense than others, especially for ordinal variables.)

# Choosing the Appropriate Proportion

Does there appear to be a relationship between gender and whether the student is looking for a spouse in college?

gender	looking for spouse			
	No	Yes	Total	
Female	86	51	137	
Male	52	18	70	
Total	138	69	207	

To answer this question we examine the row proportions:

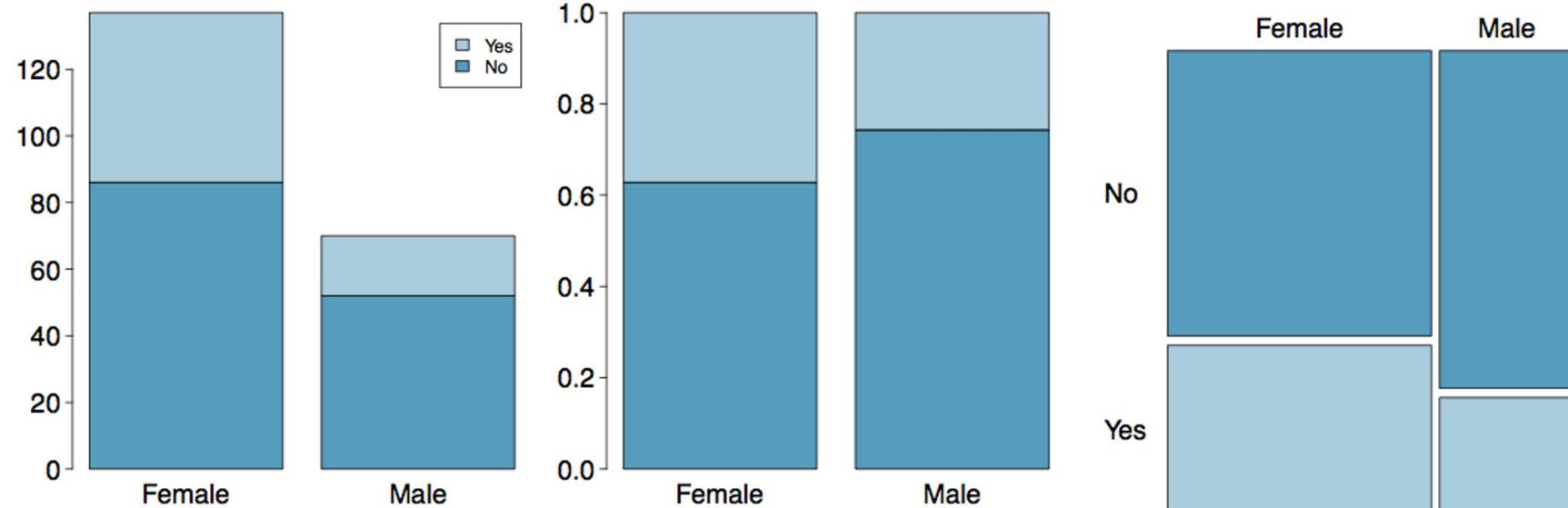
- % Females looking for a spouse:  $51 / 137 \sim 0.37$
- % Males looking for a spouse:  $18 / 70 \sim 0.26$

# Bar plots with two variables

- *Stacked bar plot*: Graphical display of contingency table information, for counts.
- *Side-by-side bar plot*: Displays the same information by placing bars next to, instead of on top of, each other.
- *Standardized stacked bar plot*: Graphical display of contingency table information, for proportions.

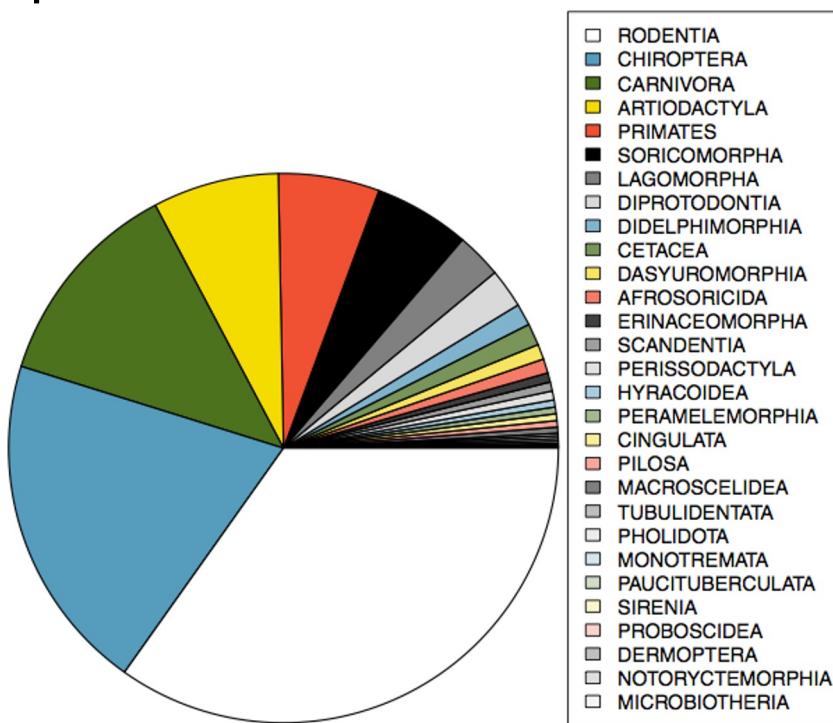
# Segmented Bar and Mosaic Plots

What are the differences between the three visualizations shown below?



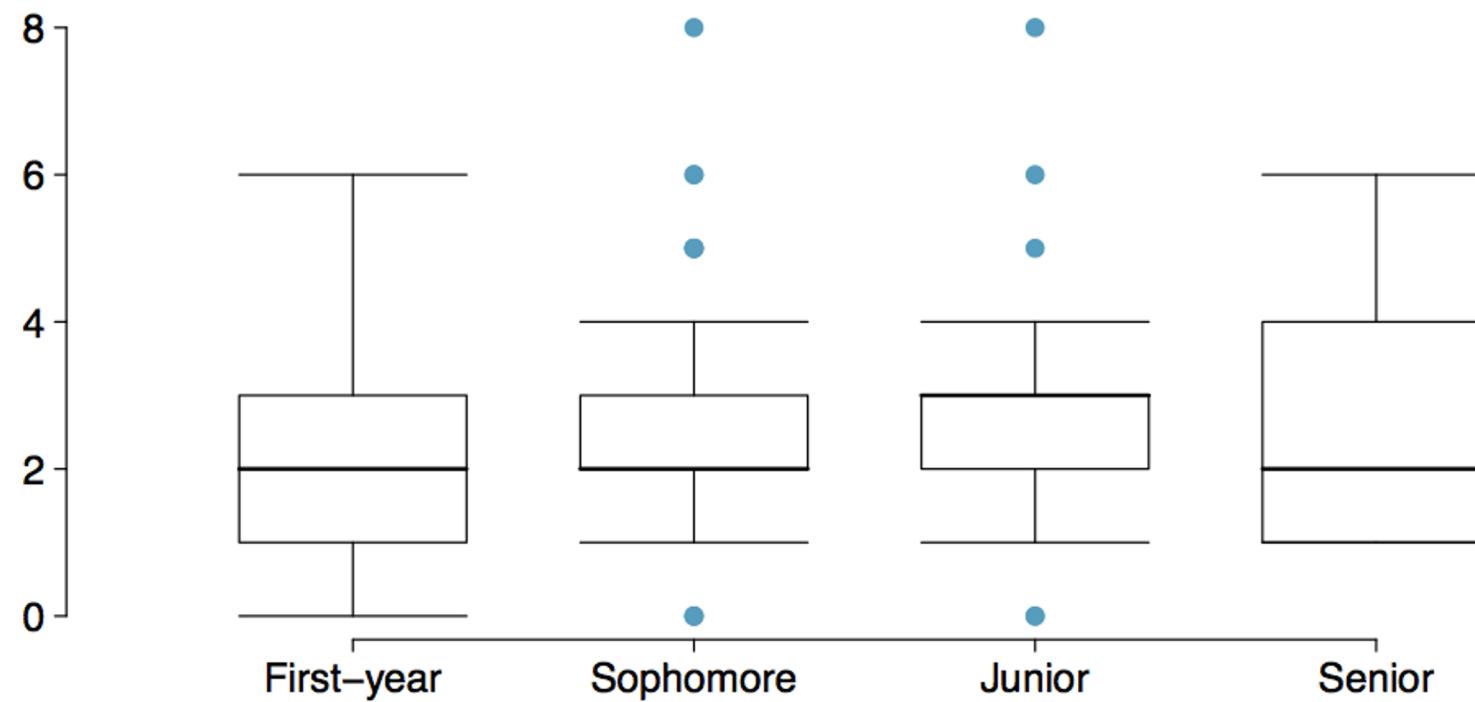
# Pie Charts

Can you tell which order encompasses the lowest percentage of mammal species?



# Comparing Numerical Data Across Groups

Does there appear to be a relationship between class year and number of clubs students are in?

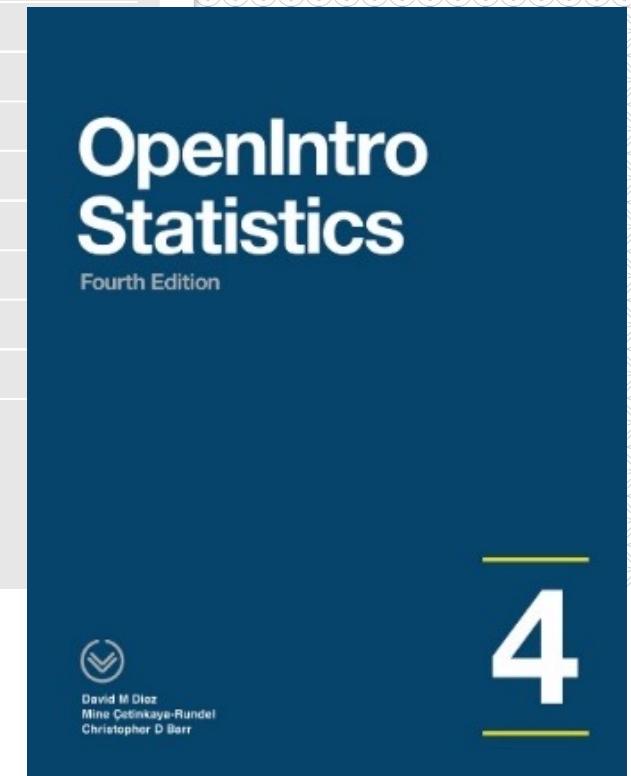


# Reference

openintro.org/stat/teachers.php



The screenshot shows the OpenIntro Statistics website's "TEACHERS" section. On the left sidebar, there are links for Register, Essentials, Lecture Slides (which is highlighted in yellow), Sample Exams, Sample Syllabuses, Learning Objectives, and Educational Software. The main content area has a header "LECTURE SLIDES" and a message stating "The slides have been updated for the 4th Edition. All LaTeX and Google Slides versions are now available." Below this, a numbered list of chapters is provided: 1 - Introduction to Data, 2 - Summarizing Data, 3 - Probability, 4 - Distributions of Random Variables, 5 - Foundations for Inference, 6 - Inference for Categorical Data, 7 - Inference for Numerical Data, 8 - Intro to Linear Regression, and 9 - Multiple and Logistic Regression.



The Slides developed by Mine Cetinkaya-Rundel of OpenIntro. Some slides adapted from Dr. Baharak course

<https://www.openintro.org/stat/teachers.php>