



The background of the slide is a blue-toned financial candlestick chart. It features several white and light blue candlesticks representing price movements. Overlaid on the chart are two main trend lines: a purple line curving upwards from the bottom left and a blue line curving downwards from the top right. A series of horizontal lines representing Fibonacci retracement levels are drawn across the chart. Specific values are highlighted in green callout boxes: '116.71' at the top left, '86.72' at the bottom left, and '99.19' near the center. Retraction percentages and their corresponding values are also labeled: '38.2%: 119.29', '51.25%: 108.98', and '61.6%: 99.19'. The title 'Intro to Data and Descriptive Statistics' is centered in a white box with a black border.

Intro to Data and Descriptive Statistics

Today:

- ▶ What are different types of data?
- ▶ What are descriptive statistics and what job do they do?
- ▶ Which descriptive statistics are appropriate for which type of data?

Data Types: the 'Levels of Measurement'

▶ **Nominal;**

- ▶ Qualitative classification of different objects by names – measures membership;
- ▶ Examples: Gender, nationality, zip code, eye color, error code;

▶ **Ordinal;**

- ▶ Categories with a natural ordering, but no well-defined scale – measures rank;
- ▶ Examples: Party membership, polling agreement (Likert) scales, ed level, class;

▶ **Interval;**

- ▶ Difference btwn units on scale is constant, but no zero point – measures exact difference;
- ▶ Examples: Time of day, date, temperature (F or C), test scores, IQ;

▶ **Ratio;**

- ▶ Difference btwn units on scale is constant/has a zero point – measures exact difference +;
- ▶ Examples: Height and weight, earnings, military spending, tax rate, temperature (K).

What are descriptive statistics?!

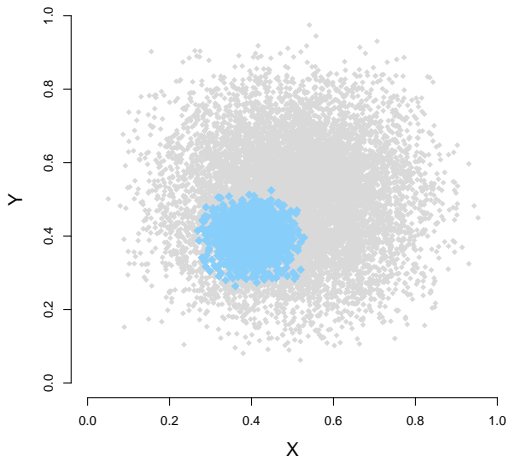
- ▶ In general, two kinds of statistics:
 - ▶ **Descriptive Statistics** – what we'll talk about today;
 - ▶ **Inferential Statistics** – what we'll spend much of the rest of the semester on;
- ▶ Typically, descriptive statistics are always reported even if main focus is on something more sophisticated.

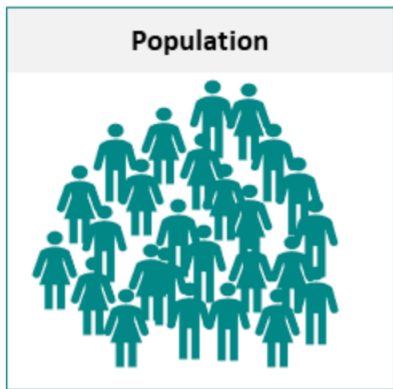
What are descriptive statistics: key terms

- ▶ **Population:** a 'complete' group of N objects, items, entities, or events of interest – e.g. all adults living in the US;
- ▶ **Sample:** a selected subset of n individuals from a population – e.g. 5,000 US adults appearing in a poll;
- ▶ **Summary Statistic:** a summary of the information in a set of observations – e.g. mean, median, mode, etc.;
- ▶ **Parametric:** derived from a probability distribution – e.g. a z-score is related to a normal distribution;
- ▶ **Nonparametric:** NOT derived from a probability distribution – e.g. descriptive statistics, histograms, etc.;
- ▶ **Univariate:** dealing with a single variable;
- ▶ **Multivariate:** dealing with relationships between several variables;

What job do they do?

- ▶ The main job of descriptive statistics is to **summarize** the information in a **sample**;
 - ▶ ...describe the data in the sample;
 - ▶ ...assess data quality (e.g. variation, correlation btwn variables, etc.);
 - ▶ ...support later inferential analysis;
- ▶ The main job of inferential statistics is to **learn** about the **population** that the sample comes from.





Measures of Central Tendency (for dataset x : 0, 1, 4, 4, 5, 5, 7, 7, 7, 9, 9)

- **Mode** – the **sample** mode is written as \bar{x}_{mode} and is the element that occurs most often in the sample. In our example $\bar{x}_{mode} = 7$.

Measures of Central Tendency (for dataset x : 0, 1, 4, 4, 5, 5, 7, 7, 7, 9, 9)

- **Median** – the **sample** median is written as \bar{x}_{med} :

$$\bar{x}_{med} = \begin{cases} x_{(n+1)/2} & \text{if } n \text{ is odd} \\ x_{(n/2)} + x_{(n/2)+1} & \text{if } n \text{ is even} \end{cases} \implies \bar{x}_{med} = x_{(11+1)/2} = x_6 = 5.$$

Measures of Central Tendency (for dataset x : 0, 1, 4, 4, 5, 5, 7, 7, 7, 9, 9)

- **Arithmetic mean** – the **sample** mean is written as \bar{x} :

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \Rightarrow \frac{0 + 1 + 4 + 4 + 5 + 5 + 7 + 7 + 7 + 9 + 9}{11} \approx 5.273.$$

Measures of Central Tendency (for dataset x : 0, 1, 4, 4, 5, 5, 7, 7, 7, 9, 9)

- **Midrange** – the **sample** midrange is written as \bar{x}_{mid} :

$$\bar{x}_{mid} = \frac{\max\{x\} + \min\{x\}}{2} \implies \bar{x}_{mid} = \frac{9 + 0}{2} = 4.5.$$

Measures of Variability (for dataset x : 0, 1, 4, 4, 5, 5, 7, 7, 7, 9, 9)

- **Range** – written as σ_{\max} , the distance between the min and max:

$$\sigma_{\max} = \max\{x\} - \min\{x\} \implies 9 - 0 = 9.$$

Measures of Variability (for dataset x : 0, 1, 4, 4, 5, 5, 7, 7, 7, 9, 9)

- **Variation ratio** – written as σ_{vr} , the proportion of cases NOT in the modal category:

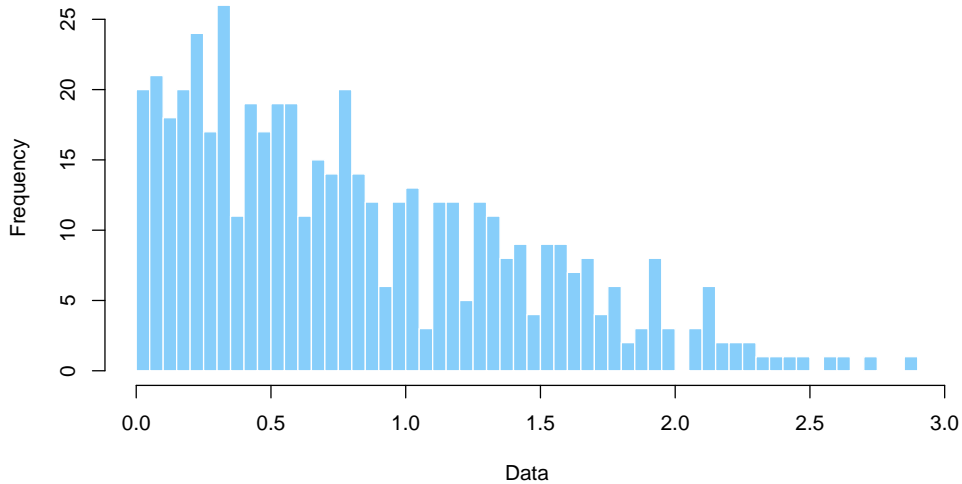
$$\sigma_{vr} = 1 - \frac{f_m}{n} \implies 1 - \frac{3}{11} \approx 0.727, \text{ where } f_m = \# \text{ of cases IN the modal category.}$$

Measures of Variability (for dataset x : 0, 1, 4, 4, 5, 5, 7, 7, 7, 9, 9)

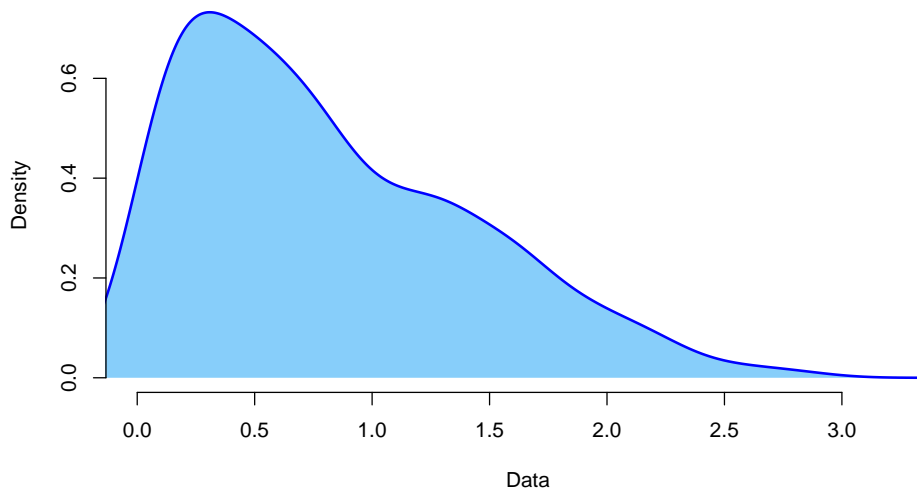
- **Standard deviation/Variance** – written as σ , the sum of squared distance from mean:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\frac{(0 - 5.273)^2 + (1 - 5.273)^2 + \dots + (9 - 5.273)^2}{11}} \approx 2.799.$$

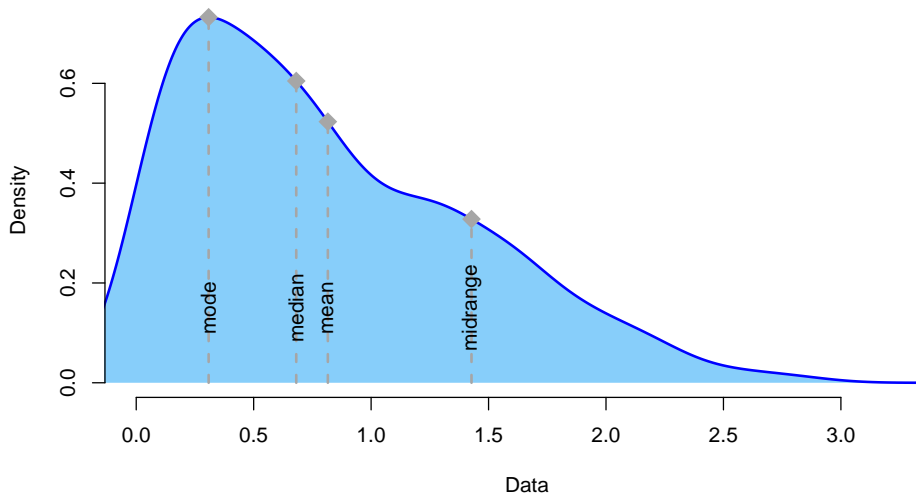
Histogram of x



Density of x



Density of x



Data Types: the 'Levels of Measurement'

► **Nominal;**

- Qualitative classification of different objects by names – measures membership;
- Examples: Gender, nationality, zip code, eye color, error code;
- Appropriate: equality, mode, Variation ratio;

► **Ordinal;**

- Categories with a natural ordering, but no well-defined scale – measures rank;
- Examples: Party membership, polling agreement (Likert) scales, ed level, class;
- Appropriate: above plus $>$ and $<$, median, range;

► **Interval;**

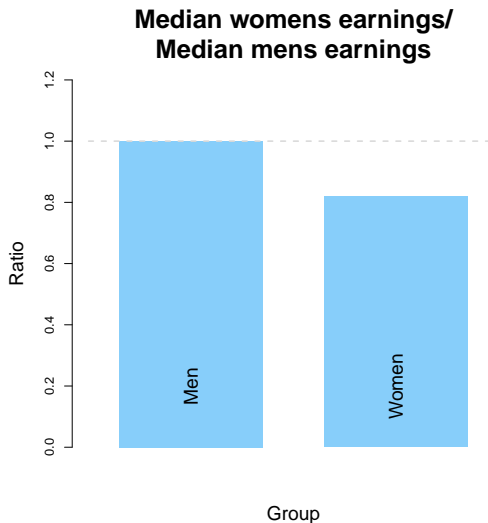
- Difference btwn units on scale is constant, but no zero point – measures exact difference;
- Examples: Time of day, date, temperature (F or C), test scores, IQ;
- Appropriate: above plus $+$ and $-$, mean, standard deviation;

► **Ratio;**

- Difference btwn units on scale is constant/has a zero point – measures exact difference $+$;
- Examples: Height and weight, earnings, military spending, tax rate, temperature (K).
- Appropriate: above plus $*$ and $/$

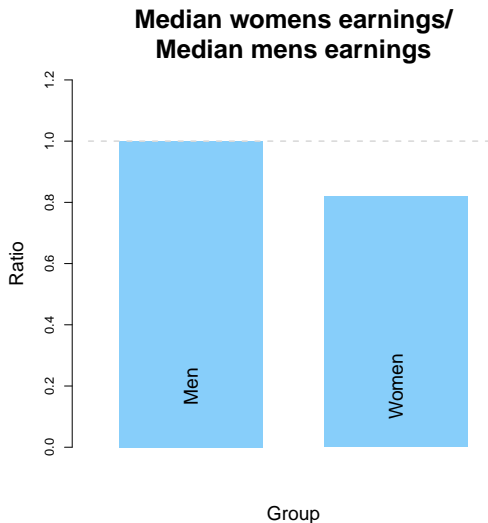
Example: Gender Wage Gap

- ▶ Lots of progress – more women in labor market with higher education than ever;
- ▶ Refers to the earnings difference between women and men:
 - ▶ Women consistently earn less than men in US;
 - ▶ But how to measure just how much less?
 - ▶ ...and what drives the gap???
- ▶ Simple descriptive statistics:
 - ▶ Compute median annual earnings for women and men working full time;
 - ▶ Take the ratio.



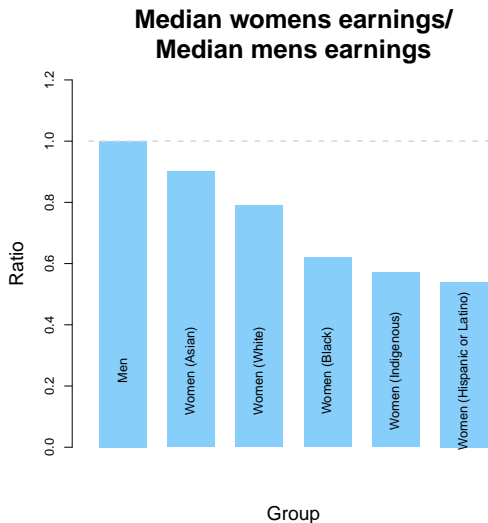
Example: Gender Wage Gap

- ▶ Lots of progress – more women in labor market with higher education than ever;
- ▶ Refers to the earnings difference between women and men:
 - ▶ Women consistently earn less than men in US;
 - ▶ But how to measure just how much less?
 - ▶ ...and what drives the gap???
- ▶ Simple descriptive statistics:
 - ▶ Compute median annual earnings for women and men working full time;
 - ▶ Take the ratio.
 - ▶ Is this enough?



Example: Gender Wage Gap

- ▶ Lots of progress – more women in labor market with higher education than ever;
- ▶ Refers to the earnings difference between women and men:
 - ▶ Women consistently earn less than men in US;
 - ▶ But how to measure just how much less?
 - ▶ ...and what drives the gap???
- ▶ Simple descriptive statistics:
 - ▶ Compute median annual earnings for women and men working full time;
 - ▶ Take the ratio.
 - ▶ Is this enough? NO!!!



Why should we care?

Using the right descriptive statistics for your data is a good and EASY first pass.