

# CSCI 3022-002 Intro to Data Science

## Exploratory Data Analysis

### Opening Zoom poll:

How do you feel about **SUMS**?  $\Sigma$

# Announcements and To-Dos

## Announcements:

1. nb day Friday! Hype!

## Last time we learned:

1. How to Zoom poll.

## To do:

1. Check out nb00 and numpy/pandas tutorials to ensure you have everything working and are getting familiar with the language used in class!
2. Read the Syllabus! It's on canvas!

# Populations and Samples

predict weight  
function of height

Statisticians hope to learn about some characteristic/variable in a population. But we often can't see the whole population; so, we investigate a sample.

**Definition:** *Population* (abstraction: every person ever in the world)

A *population* is a collection of units (units can be people, widgets, servings of food, kittens, songs, Tweets, etc.)

**Definition:** *Sample* (data).

A *sample* is a subset of the population.

**Definition:** Variable of Interest (Vol)

A *characteristic/variable of interest* is something to be measured for each unit.

# Populations and Samples

Statisticians hope to learn about some characteristic/variable in a population. But we often can't see the whole population; so, we investigate a sample.

**Example:** Suppose CU wants to determine the happiness of CS students by a survey.

1 Population

CS students

past, present, future

2 Sample → some small # of past/present students

3 Var → Happiness (score?)

# Populations and Samples

Statisticians hope to learn about some characteristic/variable in a population. But we often can't see the whole population; so, we investigate a sample.

**Example:** Suppose CU wants to determine the happiness of CS students by a survey.

## 1 *Population*

1a CSCI students, present and future

## 2 *Sample*

2a 1 in 5 current students polled, less than half respond

## 3 *Vol*

3a Happiness (a Likert scale?)

# Types of Samples

*Candidates: current + CS students*

- ▶ Simple random sample: randomly select people from sample frame *Each and every person is equally likely to have been selected.*
- ▶ Systematic sample: order the sample frame. Choose integer  $k$ . Sample every  $k$ th unit in the sample frame.
- ▶ Census sample: sample literally everyone/everything in the population *(in the frame)*
- ▶ Stratified sample: if you have a heterogeneous population that can be broken up into homogeneous groups, randomly sample from each group proportionate to their prevalence in the population.

## Inference and Generalizability

Statisticians learn about a characteristic in a population by studying a **sample**.

A major component of this course is to figure out how they make the jump from sample to population— Statistical Inference!

Statistical inference is can be informally thought of as *the study of missing information*.

## Inference and Generalizability

Statisticians learn about a characteristic in a population by studying a **sample**.

A major component of this course is to figure out how they make the jump from sample to population— Statistical Inference!

Statistical inference is can be informally thought of as *the study of missing information*.





# Exploratory Data Analysis

Before we learn about inference, we're first going to learn how to explore data. This is helpful for summarizing, recognizing patterns, etc.

There are two main types of explorations: numerical and *graphical*.

## Numerical Summaries

The calculation and interpretation of certain summarizing numbers can help us gain a better understanding of the data.

These sample numerical summaries are called **sample statistics**.

# Measures of Centrality

Summarizing the “center” of the sample data is a popular and important characteristic of a set of numbers. The goal here is to capture something like the “typical” unit with respect to the Vol.

The three most popular measures for centrality

1. The mean *-average (arithmetic average)*
2. The median
3. The mode

# The Sample Mean

$X = \text{np.array} [ \quad ]$

$\bar{x} =$

$X$  is a list of  $\#s$ , our sample

## Definition: Mean

For a given set of  $n$  numbers (observations)  $X_1, X_2, \dots, X_n$ , the sample mean or arithmetic average is

$$\bar{X} = \frac{X_1 + X_2 + X_3 + \dots + X_n}{n} = \frac{\sum_{i=1}^n X_i}{n}$$

### 1. Advantages:

① Cost  $O(n)$

② Includes outliers  $\rightarrow$  uses all the data for better, or for worse.

### 2. Disadvantages:

① Not Robust: one point can ruin  $\bar{X}$ .

② Includes outliers

# The Sample Mean

**Definition:** *Mean*

For a given set of  $n$  numbers (observations)  $X_1, X_2, \dots, X_n$ , the sample *mean* or *arithmetic average* is

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

1. Advantages:
2. Disadvantages:

# The Sample Mean

**Definition:** *Mean*

For a given set of  $n$  numbers (observations)  $X_1, X_2, \dots, X_n$ , the sample *mean* or *arithmetic average* is

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

1. Advantages:  
“Easy” to calculate; uses all data;
2. Disadvantages:  
Outliers can matter quite a bit!

# The Sample Median

## Definition: Median

For a given set of  $n$  numbers (observations)  $X_1, X_2, \dots, X_n$ , the sample *median* is the middle observation when ordered smallest to largest.

More formally, for data ordered smallest to largest  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ :

$$\tilde{x} = \begin{cases} X_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ \text{Average of } X_{(\frac{n}{2})}, X_{(\frac{n}{2}+1)} & \text{if } n \text{ even} \end{cases}$$

$X_n$  # $n$   
 $X_{(n)}$   $n$ th in order, when sorted

1. Advantages
2. Disadvantages

# The Sample Median

## Definition: Median

For a given set of  $n$  numbers (observations)  $X_1, X_2, \dots, X_n$ , the sample *median* is the middle observation when ordered smallest to largest.

More formally, for data *ordered* smallest to largest  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ :

$$\tilde{x} = \begin{cases} X_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ \text{Average of } X_{(\frac{n}{2})}, X_{(\frac{n}{2}+1)} & \text{if } n \text{ even} \end{cases}$$

### 1. Advantages

Not using all data makes it less impacted by single observations

### 2. Disadvantages

Not using all data makes it less impacted by single observations

*(thrown out info)*

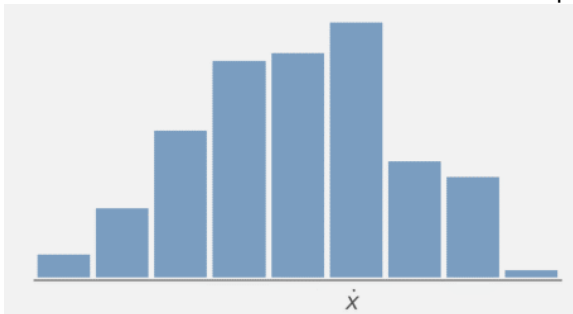


# The Sample Mode

**Definition:** *Mode*

The sample *mode* is the value that occurs the most often in the sample.

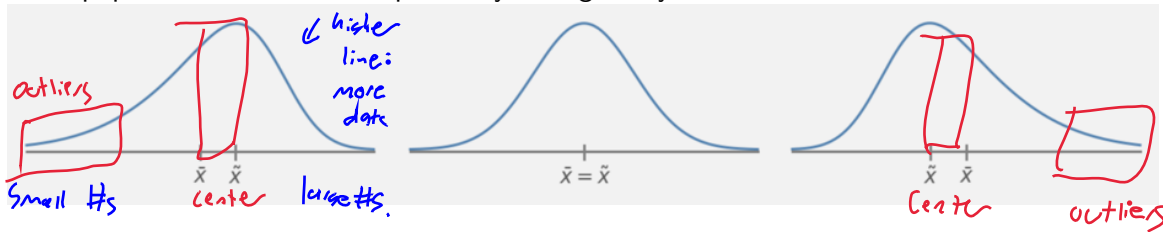
Count  
data!



# Skewness: The Mean Versus the Median

The population mean and median will generally not be equal.  
If the population distribution is positively or negatively skewed...

→ effect of outliers



Mean < Median  
"Left skew"

Mean  $\approx$  Median  
"Symmetric"

Mean > Median  
"Right skew"

precipitation

## Quartiles

**Definition:** *Quartiles* Divide the data into 4 equal parts

## Quartiles

**Definition:** *Quartiles* Divide the data into 4 equal parts

**Example:** Calculations of the median and quartiles:

Calculate the sample median and quartiles of the data:

36, 15, 39, 41, 40, 42, 47, 49, 7, 6



Quartiles:  $\begin{pmatrix} 15 \\ 39.5 \\ 42 \end{pmatrix}$

*Quartiles and Percentiles* are generalizations of quartiles.

## Quartiles

**Definition:** *Quartiles* Divide the data into 4 equal parts

**Example:** Calculations of the median and quartiles:

Calculate the sample median and quartiles of the data:

36, 15, 39, 41, 40, 42, 47, 49, 7, 6

First, sort the data: 6, 7, 15, 36, 39, 40, 41, 42, 47, 49

*Quantiles* and *Percentiles* are generalizations of quartiles.

## Quartiles

**Definition:** *Quartiles* Divide the data into 4 equal parts

**Example:** Calculations of the median and quartiles:

Calculate the sample median and quartiles of the data:

36, 15, 39, 41, 40, 42, 47, 49, 7, 6

First, sort the data: 6, 7, 15, 36, 39, 40, 41, 42, 47, 49

Now chop it up!: 6, 7, 15, 36, 39, *MIDDLE*, 40, 41, 42, 47, 49

*Quantiles* and *Percentiles* are generalizations of quartiles.

## Quartiles

**Definition:** *Quartiles* Divide the data into 4 equal parts

**Example:** Calculations of the median and quartiles:

Calculate the sample median and quartiles of the data:

36, 15, 39, 41, 40, 42, 47, 49, 7, 6

First, sort the data: 6, 7, 15, 36, 39, 40, 41, 42, 47, 49

Now chop it up!: 6, 7, 15, 36, 39, *MIDDLE*, 40, 41, 42, 47, 49

chopchop: 6, 7, 15, 36, 39, *MIDDLE*, 40, 41, 42, 47, 49

15 is the first quartile

39.5 is the median or second quartile

42 is the third quartile

*Quantiles* and *Percentiles* are generalizations of quartiles.

## Quartiles

There are multiple ways to define a quartile! Suppose we have a data set that contains:  
 $\vec{X} = \{1, 2, 3, 4, 5\}$ . What are its quartiles?

12 | 45



# Quartiles

w;k; quartile.

There are multiple ways to define a quartile! Suppose we have a data set that contains:  
 $\vec{X} = \{1, 2, 3, 4, 5\}$ . What are its quartiles?

1|2|3

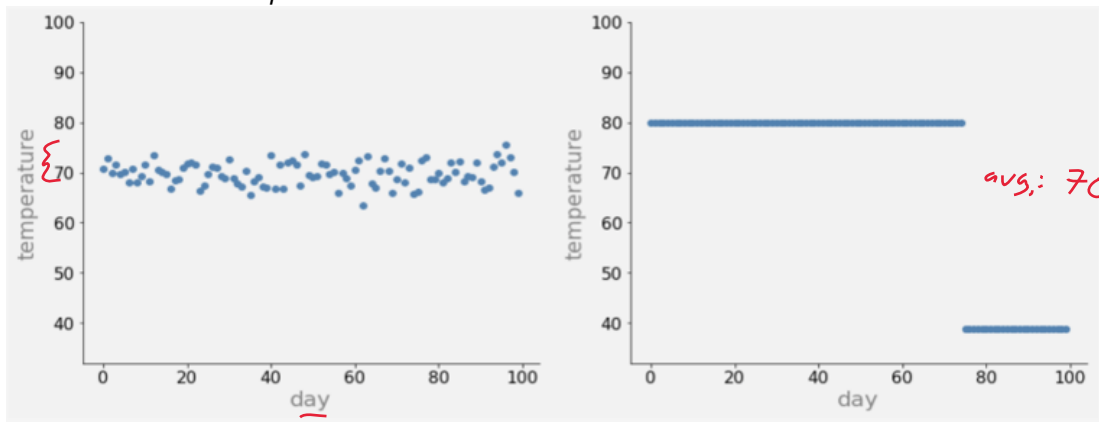
3|4|5

1. This depends on whether or not we include the median (3) in each upper and lower halves. If we do, we get 2 and 4. If we don't, we get 1.5 and 4.5.
2. There is not universal agreement between statistician *nor* software packages over which to use
3. It turns out that there's two other methods that *interpolate* the data: the median might be 1/4 of the way between two observations (1.75 and 4.25) *or* it might sit somewhere fractionally between e.g the 5/21th point and the 6/21th point. Ugh! •

## Dispersion and Spread

So far, we have learned about measuring the central tendency of data

But what about the *spread*?

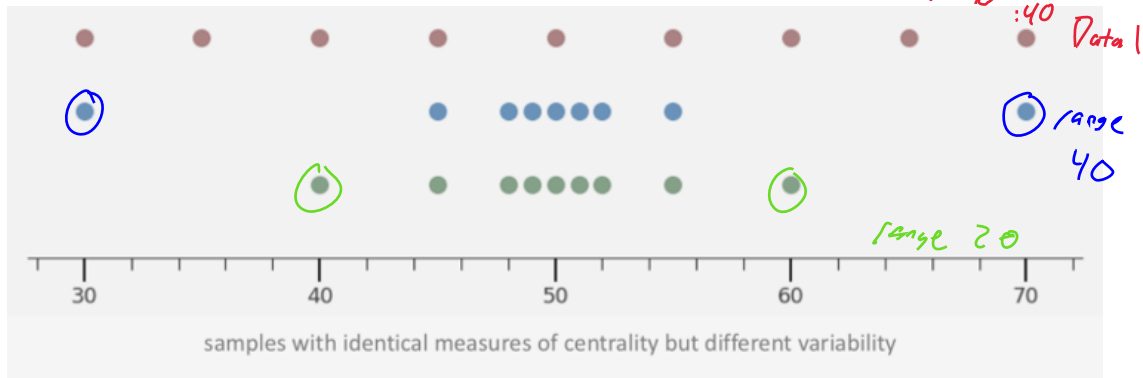


Left: San Francisco

Right: Mullensville

# The Range

Simplest measure of variability: The range.



Max - Min

## Deviation

We probably care about how far away points are from their average. “Far,” of course, is actually a math word.

- ▶ The distance between two numbers  $a$  and  $b$  is  $D = |a - b|$ .
- ▶ The distance between two points  $(a_1, a_2)$  and  $(b_1, b_2)$  is  $D = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$

We want to use the distance *from the mean*. But which type distance? Squared or not?

Dist from  $(1, 2)$  to  $(3, 7)$ :  $\sqrt{(3-1)^2 + (7-2)^2}$

## Deviation

We probably care about how far away points are from their average. “Far,” of course, is actually a math word.

- ▶ The distance between two numbers  $a$  and  $b$  is  $D = |a - b|$ .
- ▶ The distance between two points  $(a_1, a_2)$  and  $(b_1, b_2)$  is  $D = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$

We want to use the distance *from the mean*. But which type distance? Squared or not? For each datum  $X_i$ , the *deviation from the mean* of  $X_i$  is

$$|X_i - \bar{X}|$$

## Variance and Standard Deviation

### Definition: Sample Variance

The *sample variance*, denoted by  $s^2$ , is given by:

take each obs  $X_i$ , compute  $(X_i - \bar{X})$ , square and sum and average

$$s^2 = \left( \frac{n}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right) / (n-1)$$

The *sample standard deviation*, denoted by  $s$ , is the (positive) square root of the variance:

Note that  $s^2$  and  $s$  are both nonnegative. The unit for  $s$  is the same as the unit for each of the  $X_i$ .

## Variance and Standard Deviation

**Definition:** *Sample Variance*

The *sample variance*, denoted by  $s^2$ , is given by:

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

The sample *standard deviation*, denoted by  $s$ , is the (positive) square root of the variance:

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

Note that  $s^2$  and  $s$  are both nonnegative. The unit for  $s$  is the same as the unit for each of the  $X_i$ .

## Standard Deviation

**Example:** *Calculation of the SD*

Data (units in dollars): 2,4,3,5,6,4.



## Standard Deviation

**Example:** *Calculation of the SD*

Data (units in dollars): 2,4,3,5,6,4.

Since we mean business, we need the average first.

$$\bar{X} = \frac{2 + 4 + 3 + 5 + 6 + 4}{6} = \frac{24}{6} = 4$$

Now let's compute the deviations...

vectorized deviations

$$\overbrace{[(X - \bar{X})^2]}^{\text{vectorized deviations}} = [(2 - 4)^2, (4 - 4)^2, (3 - 4)^2, (5 - 4)^2, (6 - 4)^2, (4 - 4)^2]$$

and sum and “average” those!

$$s^2 = \frac{4 + 0 + 1 + 1 + 4 + 0}{5} = 2$$

## The Interquartile Range

The interquartile range is defined to be the difference between the upper and lower quartiles:

$$IQR = Q_3 - Q_1$$

It's a spread measure standardly used in *box plots*, which we introduce formally next time.

## Tukey's Five Number Summary

John Tukey, father of modern EDA, advocated summarizing data sets with 5 values:

1. Min value
2. Lower quartile
3. Median
4. Upper quartile
5. Max value

Advantages:

gives the center of the data

gives the spread of the data (range in IQR)

gives an idea of skewness (compare how far away Q1 and Q3 are from median!)

## Next Time: Visual EDA!

Collapsing our data into a few descriptive numbers is pretty valuable!

...but *summary statistics* invariably throw away a lot of detail and nuance. Maybe we should consider visualizing the data to include more information?

