

Introduction to Data Science With Probability and Statistics

Lecture 2: Exploratory Data Analysis & Summary Statistics

CSCI 3022 - Summer 2020
Sourav Chakraborty
Dept. of Computer Science
University of Colorado Boulder

What will we learn today?

- Mean
- Median
- Mode
- Quartiles
- Sample Variance
- Sample Standard Deviation
- Interquartile Range (IQR)
- 5-Number Summary



- A Modern Introduction to Probability and Statistics, section 16.1-16.3

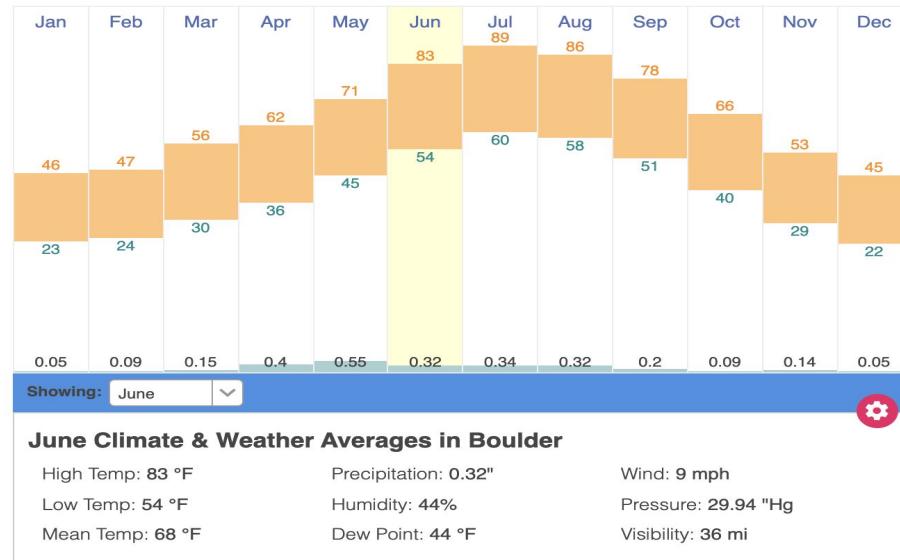
Exploratory Data Analysis

How do you summarize and describe a set of data

Annual Weather Averages Near Boulder

Averages are for Broomfield / Jeffco, which is 11 miles from Boulder.

Based on weather reports collected during 1985–2015.



* Charts taken from [TimeandData.com](https://www.timeanddate.com/weather/us/colorado/boulder)

Exploratory Data Analysis - Central Tendency

Center of a Dataset

Sample Mean: $\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$

Sample Median: The *middle element* of the dataset when it is put in ascending order.
Note: When n is even, take the average of the middle two elements.

Example: Compute the mean and median of the following set of data.

10 11 14 19 23 11 13 16 17 18 9 13 11,375

Sample Mean: $\bar{x} = \frac{10+11+14+19+23+11+13+16+17+18+9+13+11,375}{13} = 888.384615$

Sample Median: 9 10 11 11 13 13 14 16 17 18 19 23 11,375



Exploratory Data Analysis - Central Tendency

Data scientists hope to learn about some characteristic of a population, typically we cannot study the entire population so we study a **sample**.

A **population** is a collection of units. (e.g. people, songs, tweets, incomes, daily temperatures ...)

A **sample** is a subset of the population

A **characteristic/variable of interest (VOI)** is something we want to measure for each unit

Exploratory Data Analysis

Example: Suppose that the city of Boulder wants to estimate its per-household income via a phone survey. They call every 50th number on a list of Boulder phone numbers between 6PM and 8PM.

What is the **population**?

What is the **sample**?

What is the **Variable of Interest**?



Exploratory Data Analysis

The **sample frame** is the source material or device from which the sample is drawn.

Simple random sample: randomly select people from a sample frame.

Systematic sample: Order the sample frame. Choose an integer k . Sample every k^{th} unit in the sample frame.

Census sample: sample everyone/everything in the population.

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.

Exploratory Data Analysis - Central Tendency

Summarizing the “center” of the sample is a popular way to characterize data.

Measures of Central Tendency: Mean, Median, and Mode

Sample Mean: $\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{k=1}^n x_k$

Sample Median: The *middle element* of the dataset when it is put in ascending order.

Note: When n is even, take the average of the middle two elements.

Sample Mode: Most commonly occurring value in a data set.



THE MEDIAN AGE AT FIRST BIRTH IN THE US IS 25, WHICH MEANS THE TYPICAL NEW MOTHER IS NOW A 90'S KID.

Exploratory Data Analysis - Central Tendency

$$\text{Sample Mean: } \bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{k=1}^n x_k$$

Calculation:

Add up all numbers in the data set.

Divide by number of numbers.

Example: Compute the mean of the following data set.

$$\begin{array}{r} 17+10+10+18+5+23 \\ \hline \end{array}$$

$83/6 \approx 13\cdot$

Exploratory Data Analysis - Central Tendency

Sample Median: The *middle element* of the dataset when it is put in ascending order.

Note: When n is even, take the average of the middle two elements.

Calculation:

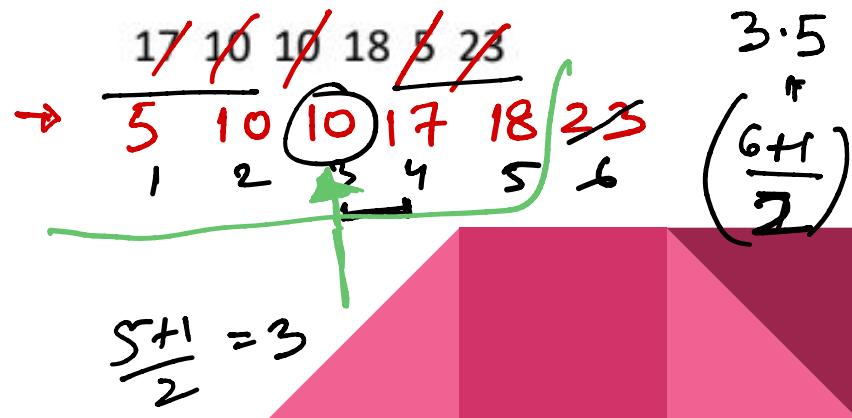
Order the n observations from smallest to largest
Include multiple instances of repeated values.

If n is odd, then $\tilde{x} = \left(\frac{n+1}{2}\right)^{th}$ ordered value.

If n is even, then $\tilde{x} = \text{average of } \left(\frac{n}{2}\right)^{th} \text{ and } \left(\frac{n}{2} + 1\right)^{th}$

$$\bar{x}, \tilde{x}$$

Example: Compute the median of the following data set.



Exploratory Data Analysis - Central Tendency

Sample Mode: Most commonly occurring value in a data set.

Calculation:

Count what number/data point occurs most often

Example: Compute the mode of the following data set.

17 10 10 18 5 23

Mode=10

1 2 3 4 5

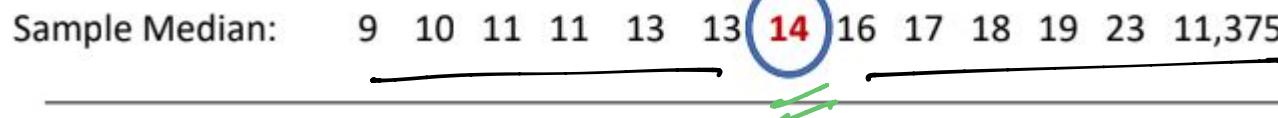
Exploratory Data Analysis - Central Tendency

Example (from earlier): Compute the mean and median of the following set of data.



9-23
14
11,375

Sample Mean: $\bar{x} = \frac{10+11+14+19+23+11+13+16+17+18+9+13+11,375}{13} = 888.384615$



Sample Mean (aka Arithmetic Average)

Advantages: Simple to compute

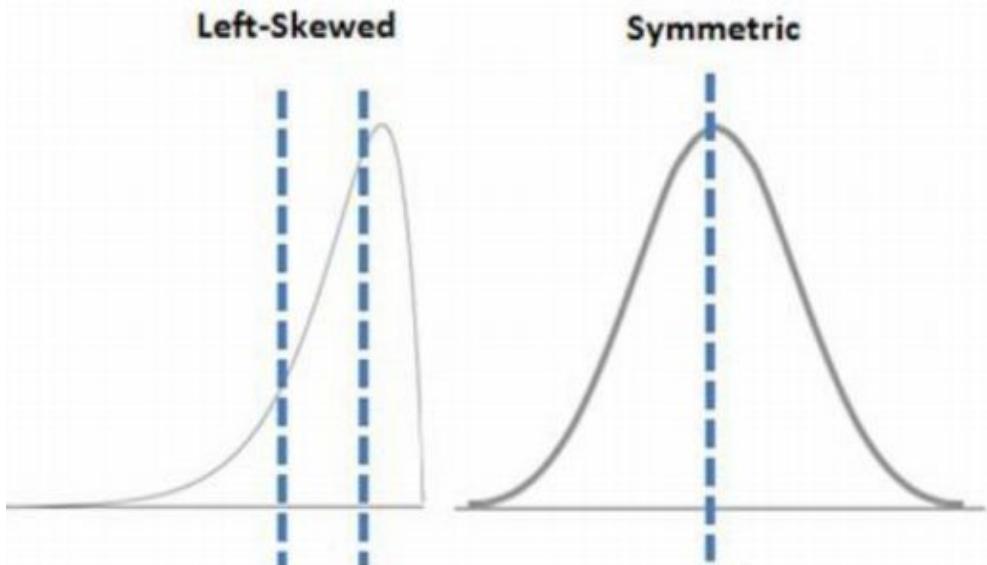
Disadvantages: Very sensitive to outliers

Sample Median

Advantages: Not as easily affected by a few outliers

Disadvantages: Doesn't depend on all entries in data set

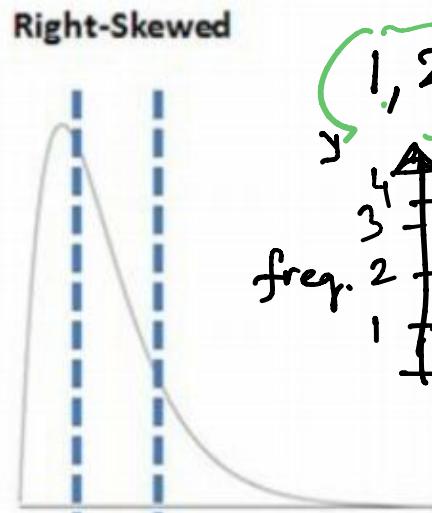
Exploratory Data Analysis - Central Tendency



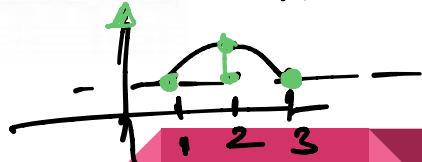
$$\bar{x} < \tilde{x}$$

$$\bar{x} = \tilde{x}$$

(1 2 2 3)



$\tilde{x} < \bar{x}$
 $\bar{x} > \tilde{x}$



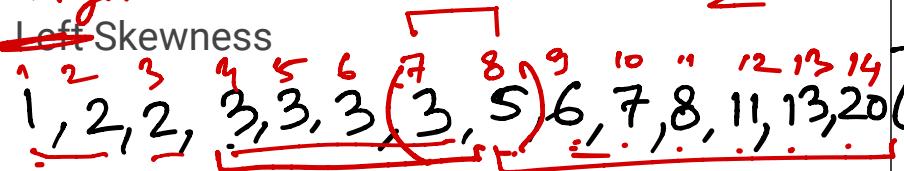
- The population mean and median will generally not be equal.

(Mean = 2 , Median = 2 , mode = 2)

Data Skewness

Right

Left Skewness



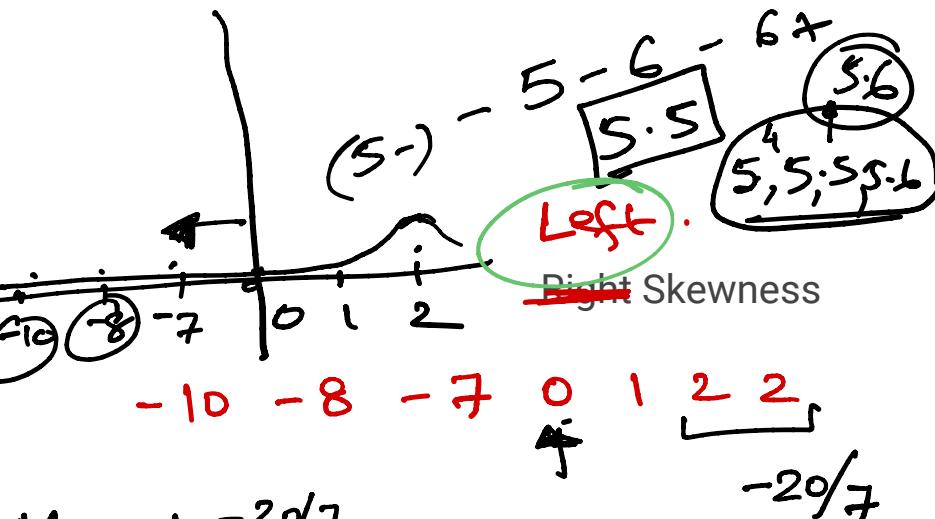
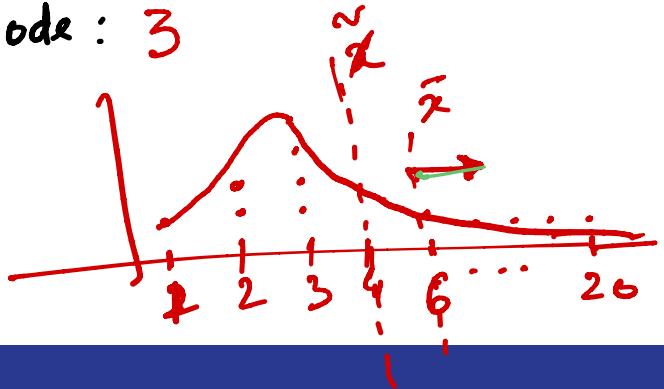
$$\frac{14+1}{2}$$

$$\frac{70+20}{14}$$

$$\text{Mean: } \approx 6 + \bar{x}$$

$$\text{Median: } 4 \quad \tilde{x}$$

$$\text{Mode: } 3$$

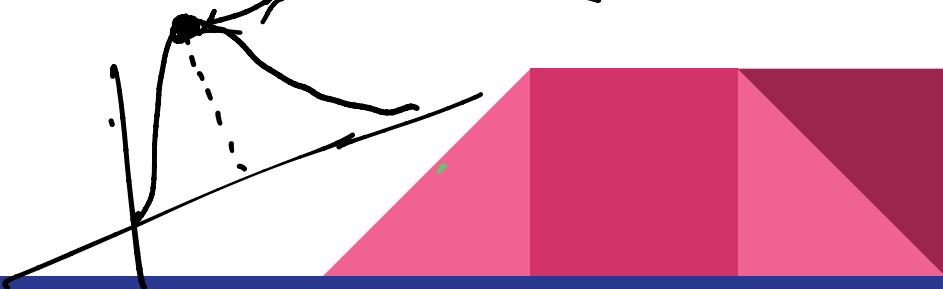


$$\text{Mean: } -20/7$$

$$\text{Median: } 0$$

$$\text{Mode: } 2$$

$$-\bar{x} < \tilde{x}$$



Exploratory Data Analysis

Other sample measures

Quartiles: Divide the data into 4 equal parts

Lower quartile (Q_1 or P_{25}) splits the lowest 25% of the data from the other 75%

Middle quartile (Q_2 or P_{50}) splits the data in half (i.e. the median) **Median**

Upper quartile (Q_3 or P_{75}) splits the highest 25% of the data from the lowest 75%

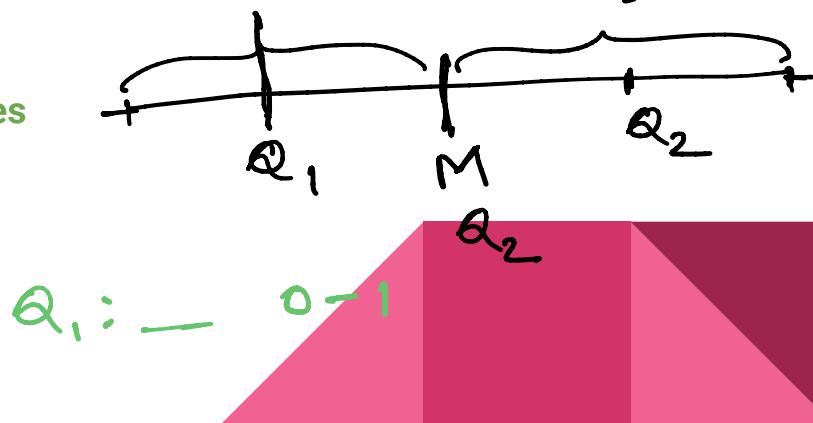
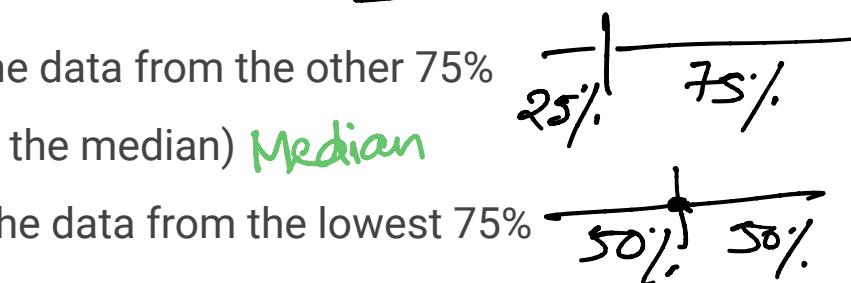
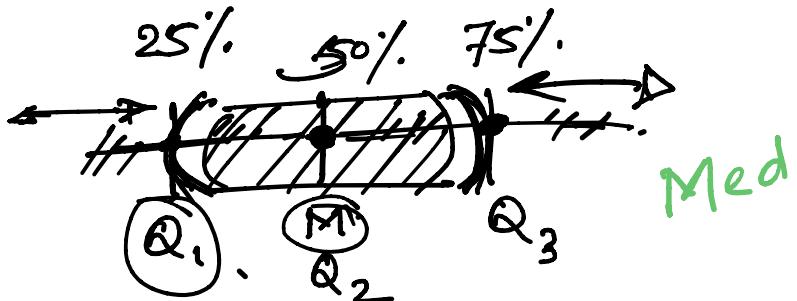
Computation:

1) Use the median to divide the ordered data set into 2 halves

- If n is odd, include the median in both halves.
- If n is even, split the data exactly in half.

2) The lower quartile is the median of the lower half.

3) The upper quartile is the median of the upper half.



Exploratory Data Analysis

Example: Compute the quartiles of the data: 6, 7, 15, 36, 39, 40, 41, 42, 43, 47, 49, 50

$$((11+1)/2) = 6$$
$$Q_1: 15$$
$$Q_3: 40.5$$

n is odd, median is a single value,
 Q_1, Q_2

n is even, median is an avg

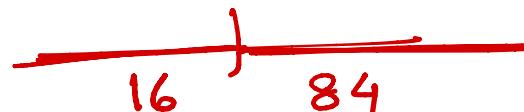
Exploratory Data Analysis

Quartiles: Divide the data into 4 equal parts

- **Lower quartile** (Q_1 or P_{25}) splits the lowest 25% of the data from the other 75%
- **Middle quartile** (Q_2 or P_{50}) splits the data in half (i.e. the median)
- **Upper quartile** (Q_3 or P_{75}) splits the highest 25% of the data from the lowest 75%

Percentiles:

- Generalization of quartiles.
- Q_1 is the 25th percentile; P_{25}
- Can also calculate general percentiles: e.g. the 16th percentile, P_{16} , splits off the lower 16% of the data.



Quantiles Generalization

$$\dots \boxed{44} \xrightarrow[45]{45; 46}$$

$$D(6), D(7) \xrightarrow{6 \cdot 5}$$

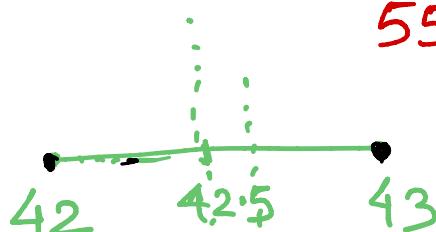
$$\frac{44}{45} \boxed{46}$$

$$\frac{n+1}{2}$$

$$= \boxed{0.5(n+1)}$$

Example: 41 41 41 41 41 42 43 43 43 58 58 $\leftarrow n$

(42, 43)



$$\text{Ex. } \tilde{x} = D\left(\frac{n+1}{2}\right)$$

$$0.55$$

$$= \underline{\underline{6}} \quad 0.5$$

$$50\%$$

$$55\%$$

$$0.55$$

$$\begin{matrix} 0.5 \\ \cdot 5 \\ \cdot 5 \\ \cdot 5 \end{matrix}$$

$$\text{id}x = \left\lfloor \frac{p \cdot (n+1)}{100} \right\rfloor \quad 0.6 * (43 - 42) = 0.6$$

$$42 + 0.6 = \underline{\underline{42.6}}$$

$$\alpha = \frac{p \cdot (n+1)}{100} - \text{id}x = 0.6$$

$$\left[D(\text{id}x) + \alpha \cdot (D(\text{id}x+1) - D(\text{id}x)) \right]$$

$$= [6 \cdot 6] = \underline{\underline{6 + 0.6}}$$

$$42, 50 \quad \boxed{6 \cdot 6}$$

$$\frac{46}{\underline{\underline{6}}}$$

$$\underline{\underline{6}}$$

Exploratory Data Analysis

$$\mathcal{D}_1 : 1, 2, 3$$

Other ways to analyze data other than central tendency?

$$\mathcal{D}_2 : -1, 2, 5$$

How about measuring the spread or variability of the data.

$$\bar{x}_1 = \bar{x}_2 = 2$$

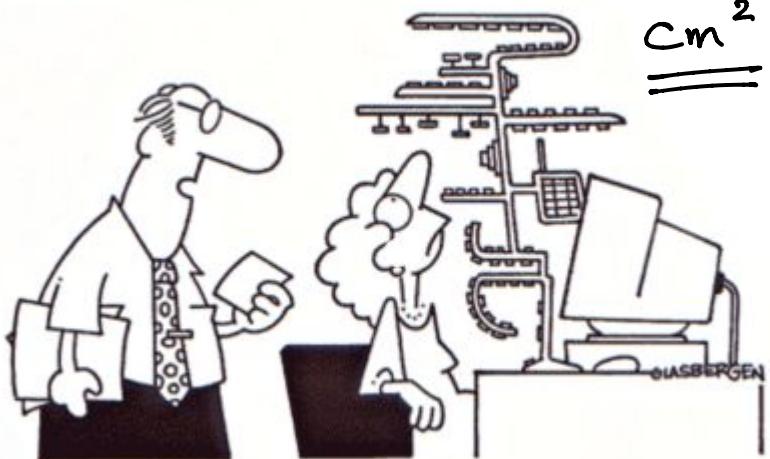
$$R_1 = 3 - 1 = 2$$

$$R_2 = 5 - (-1) = 6$$



Exploratory Data Analysis

Copyright 2002 by Randy Glasbergen.
www.glasbergen.com



"It's the new keyboard for the statistics lab. Once you learn how to use it, it will make computation of the standard deviation easier."

$$\bar{x} = 2.5$$
$$\frac{1}{n} \left[(1-2.5)^2 + (2-2.5)^2 + (3-2.5)^2 + (4-2.5)^2 \right]$$

③ (4-1)

The sample variance, denoted by s^2 , is given by

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

$$\frac{1}{n}$$

The sample standard deviation, denoted by s , is given by the square root of the variance.

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

$$\sigma = \underline{\underline{Cm}}$$

$$\checkmark$$

$$\underline{\underline{Cm}}.$$

$$\textcircled{1}$$

Exploratory Data Analysis - Variability

Example: Compute the standard deviation (SD) of the data: 2, 4, 3, 5, 6, 4

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

$$\begin{aligned}\sigma^2 &= \frac{1}{6-1} \left((2-4)^2 + (4-4)^2 + (3-4)^2 + (5-4)^2 + (6-4)^2 + (4-4)^2 \right) \\ &= \frac{1}{5} (4 + 0 + 1 + 1 + 4 + 0) = \frac{10}{5} = 2\end{aligned}$$

$$\sigma = \sqrt{\sigma^2} = \underline{\underline{\sqrt{2}}} .$$

Exploratory Data Analysis – Interquartile Range

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

$$IQR = Q_3 - Q_1$$

Assumption : Excluding the median from the computation of Q_1 , Q_3 if n is odd

- IQR gives the spread of 50% of the data

Example: Compute the IQR of the data: 6, 7, 15, 36, 39, $\overline{Q_1}$, 40, $\overline{Q_2}$, 41, 42, 43, 47, 49

$$Q_1: 15$$

$$Q_2: 40$$

$$Q_3: 43$$

$$\begin{aligned} IQR &= Q_3 - Q_1 = 43 - 15 \\ &= \boxed{28} \end{aligned}$$

Exploratory Data Analysis – 5-Number Summary

John Tukey



John Wilder Tukey

Tukey advocated summarizing data sets with 5 values:

- 1) Minimum value
- 2) Lower Quartile
- 3) Median
- 4) Upper Quartile
- 5) Maximum value

| | |
|-------------|--|
| Born | June 16, 1915 New Bedford, Massachusetts, U.S. |
| Died | July 26, 2000 (aged 85) New Brunswick, New Jersey, U.S. |

Exploratory Data Analysis – 5-Number Summary

Example: Find the 5-number summary of the data: 6, 7, 15, 36, 39, 40, 41, 42, 43, 47, 49

Min ?

Q_1 ?

Q_2 ?

Q_3 ?

Max ?

Next Time:

Box and Whisker Plots!

