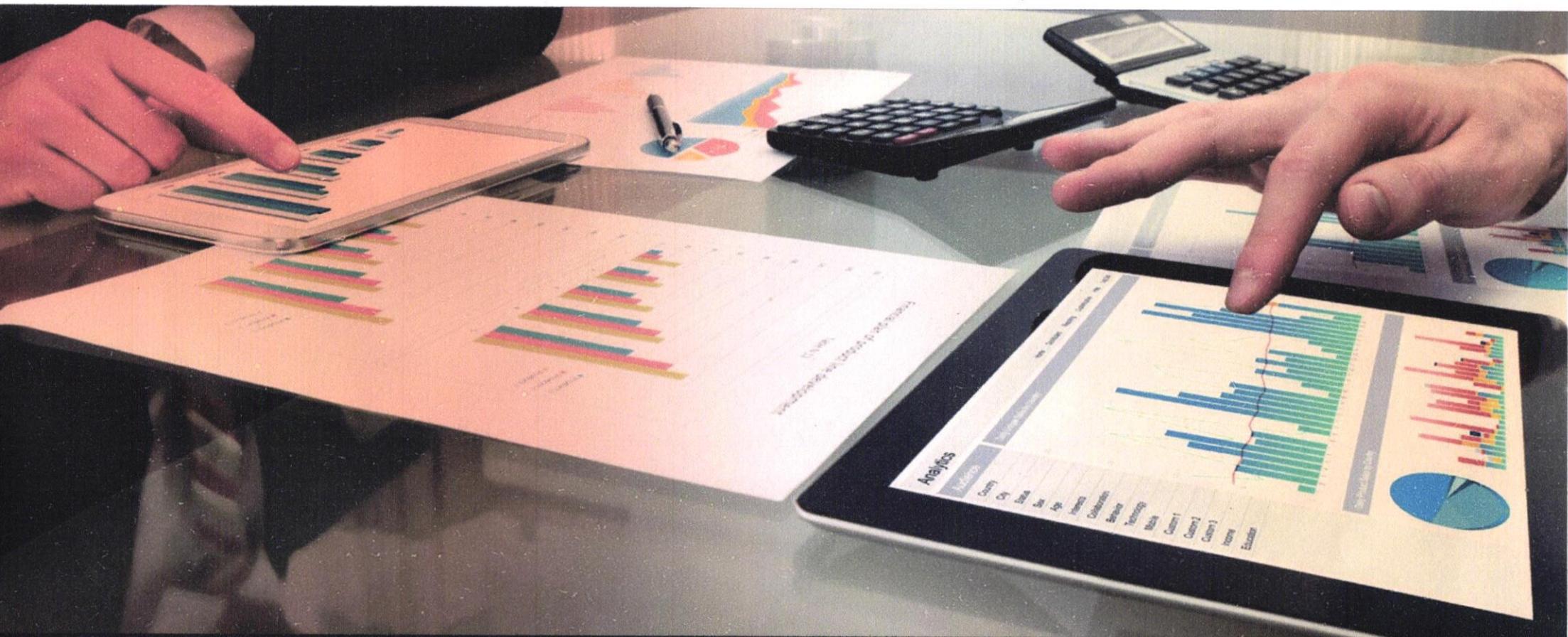


Cohort

July 2021

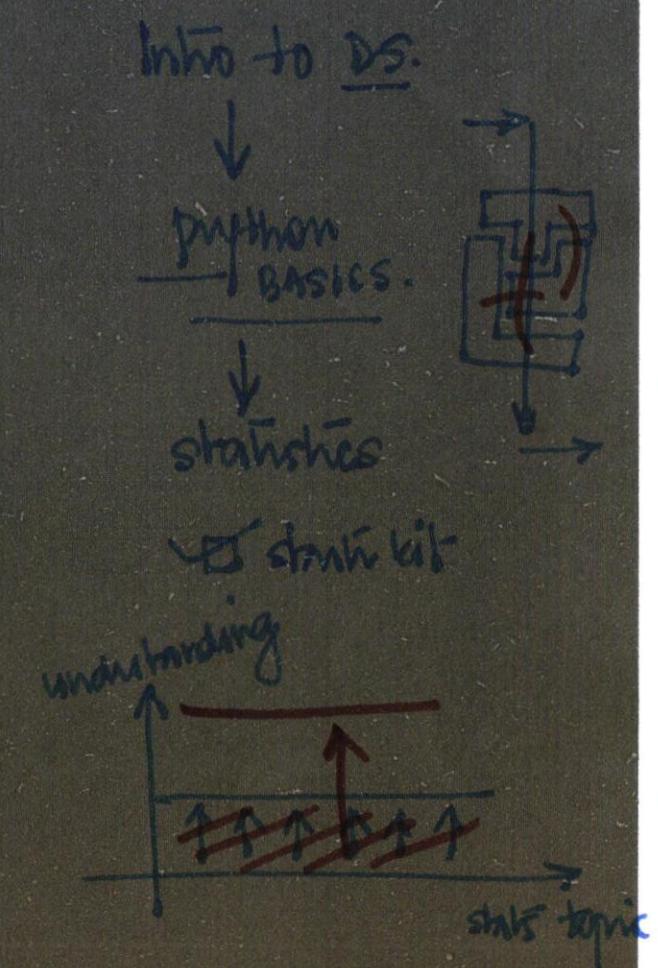


Introduction to Statistics



Agenda

- What is Statistics?
- Importance of Statistics
- Applications of Statistics
- What is Data?
- Data On Tap
- Types of Data
- Data Quality Issues
- Population vs Sample
- Types of Statistics
- Descriptive Statistics



What is Statistics?

numbers ←
engineering

analysis of
samples.

logic
structure
study of
numbers,
grammar of science.

EXIT POLLS · A prediction of
Who wins / loose?

- A form of **mathematical analysis**.
- Science of **conducting studies** to gather, review, analyze and draw conclusions from data.

Time taking
costly
representative
sample

ABC

HII

+ XYZ
HIIII

= 130
marks

130

Importance of Statistics

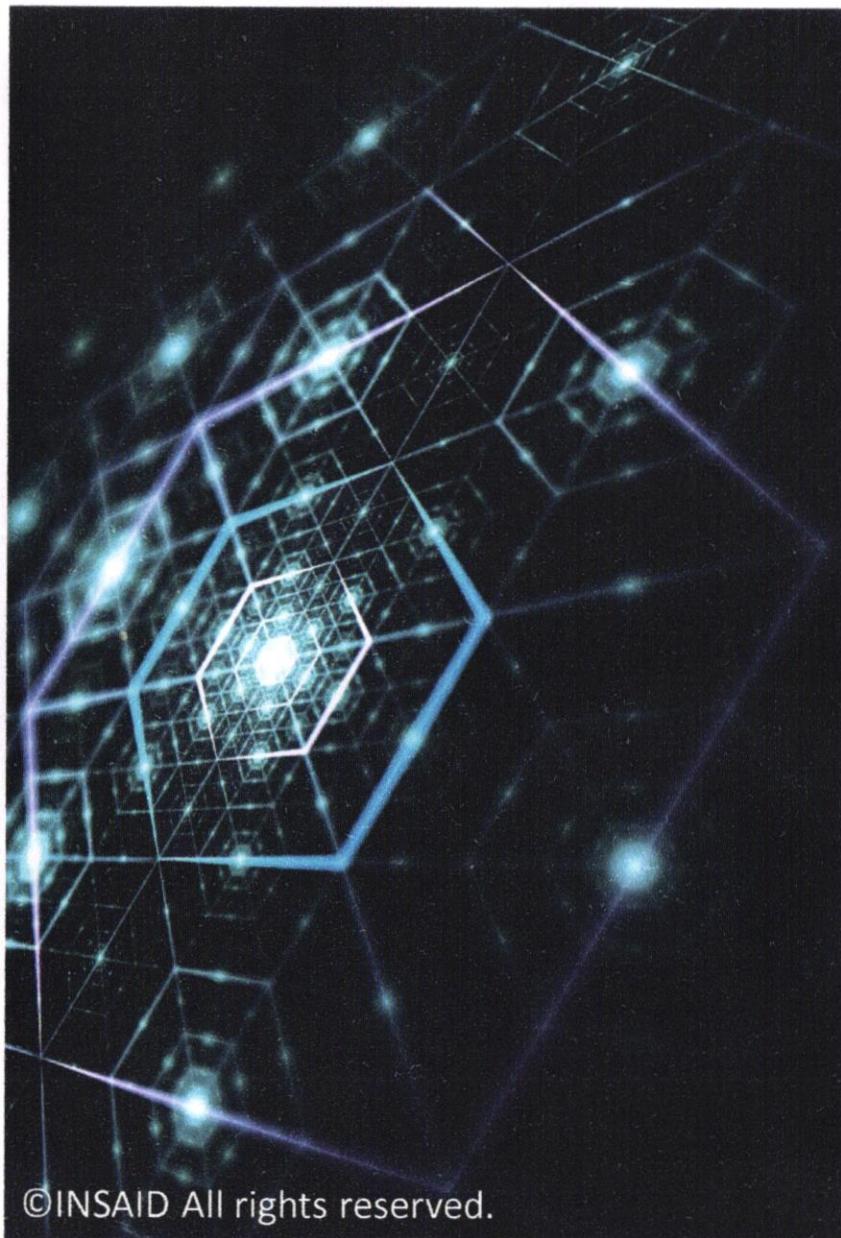
- Business Growth
- Collect the relevant data
- Knowledge Discovery
- Pattern Recognition
- Forecasting



blood pressure
120/80

Applications of Statistics

- **Medicine:** WHO research on epidemic spread
 - **Business:** Sales Forecasting
 - **Weather Forecasting:** Analysis over daily weather
 - **Stock Market:** Stock profit analysis
 - **Health:** Analysis of human genes
 - & Many More...
- # covid +s
vacination rate
Annual sales
Revenue
Temp, Hum,
Expected Rainfall
Price, monwtin,
BP, pulse rate, sugar,
ECG, CBC.



Text
====

Music
~~~~~

videos  
▶

image  
☒

Excel  
|||||

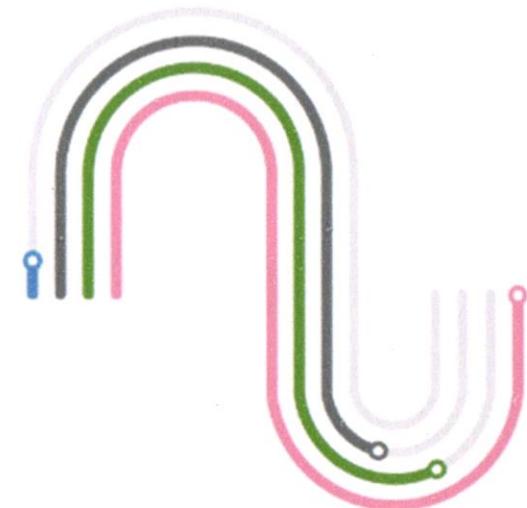
binary  
1/0

## What is Data?

anything which  
can be recorded.

Meat consumption Kg/person

|    |        |
|----|--------|
| 3  | India  |
| 79 | Brazil |
| 80 | Israel |
| 99 | US     |



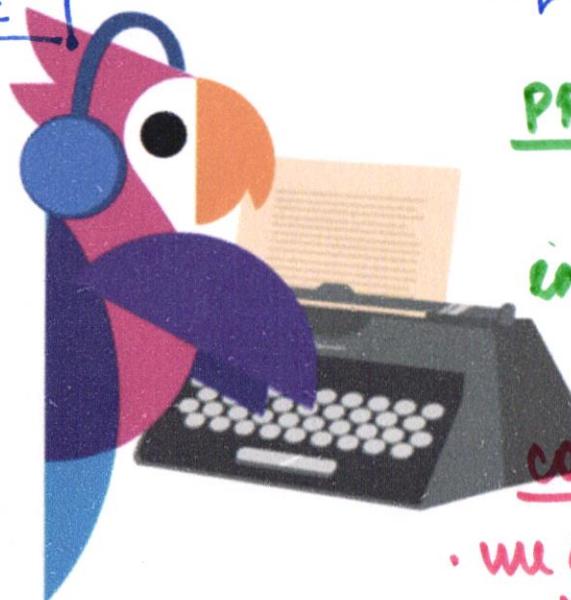
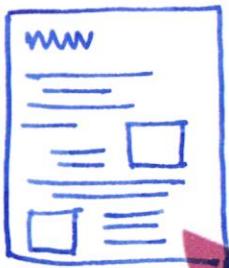
- Raw or isolated facts from which the required information is produced.
- Can be recorded & have implicit meaning.
- Example: Meat consumption Kg/person

# Data On Tap

Wire  
Internet

- Generally, data is available in two forms:
  - Primary Form
  - Secondary Form





data  
|  
**primary**

**Secondary**

→ 2<sup>nd</sup> hand data .

PRO  
cheap  
easy  
fast

CON

- we don't design the rep
- Sampling can not be dictated
- Trust issues.

→ collect ourselves from source

PRO

• clear understanding

## Primary Form

• sampling by

use

concreteness ↑

CON

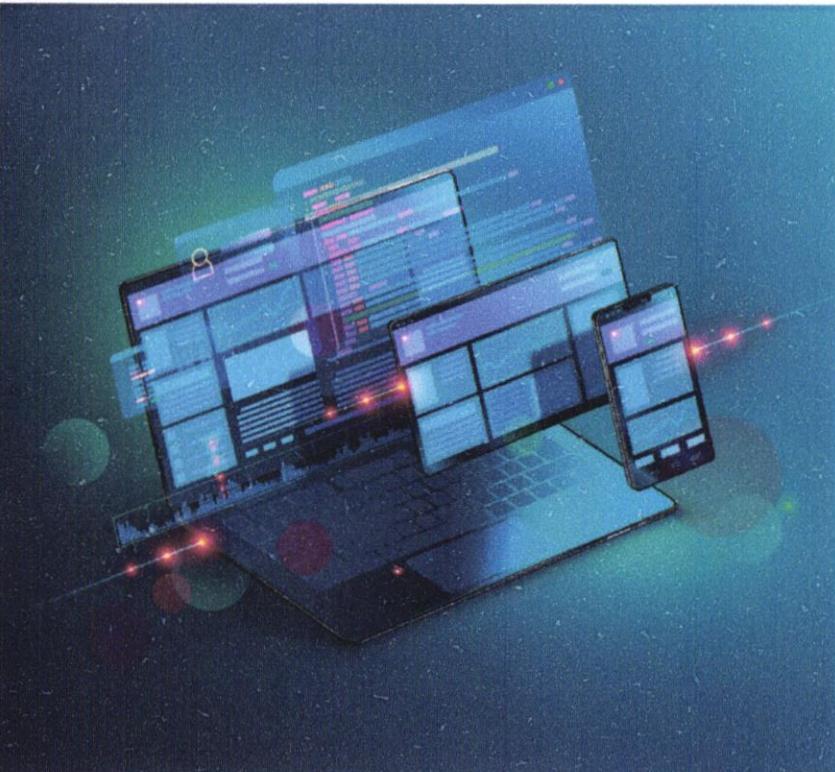
costly

Time consuming

- Collected first hand by the researcher.

- It can be acquired in the form of -

- Interviews, Observation, Questionaries, Case Studies, Group Discussions etc.

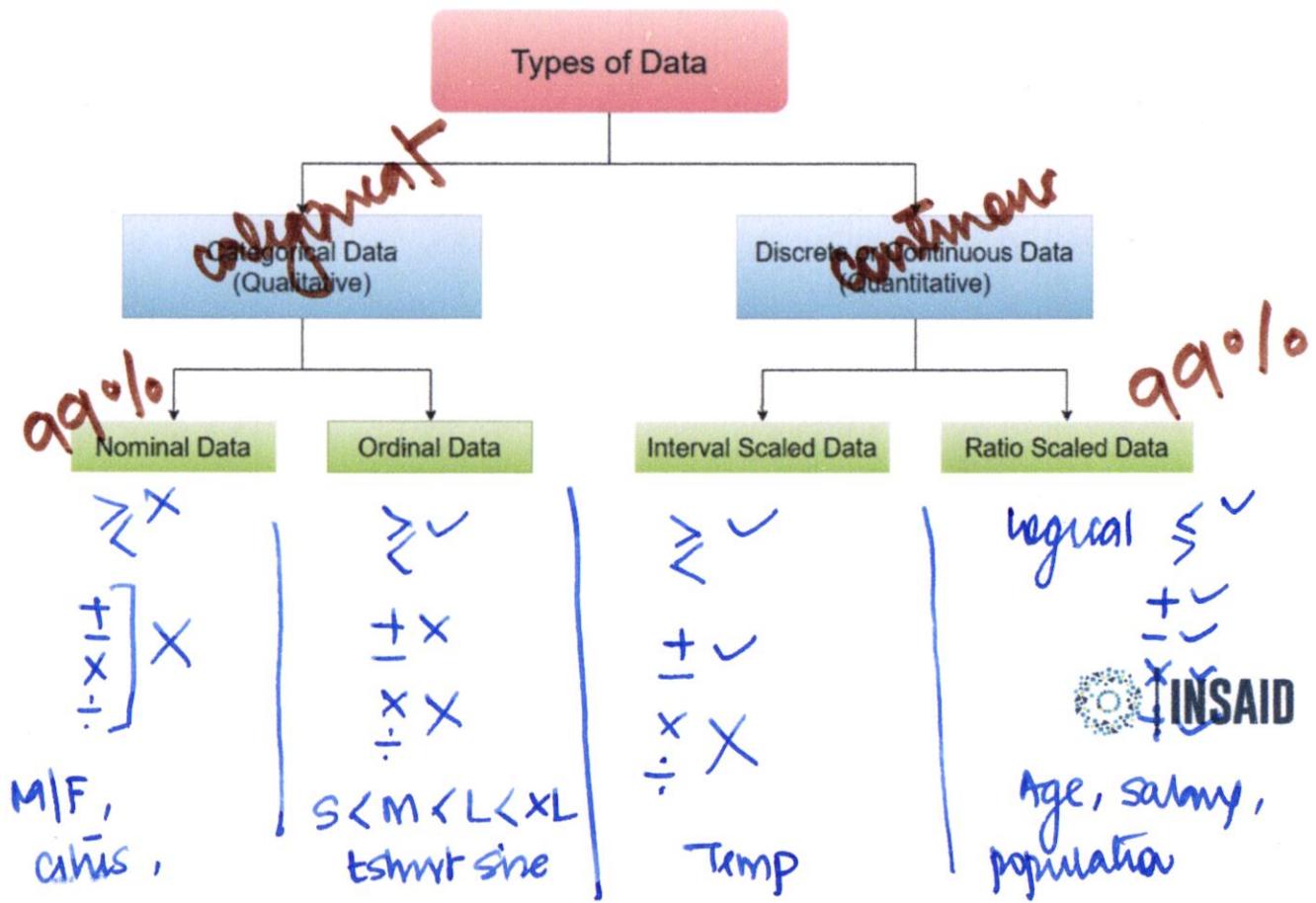


## Secondary Form

- Already collected by and readily available from other sources.
- It can be acquired in the form of -
  - Previous Research, Diaries, Letters, Web Info, Census Data etc.

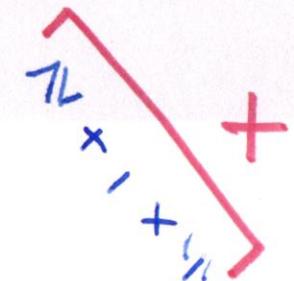


# Types of Data





## Nominal Data



- Data related to names.
- The values are symbols or names.
- Also referred as Categorical Data.
- **Example:** Hair Color

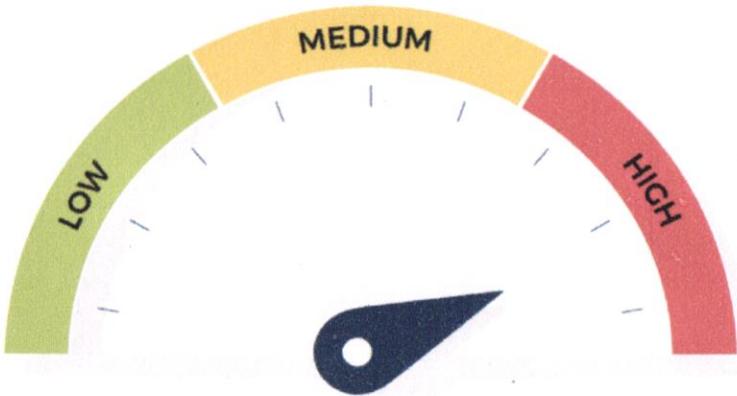
Example of  
Nominal

## Binary Data

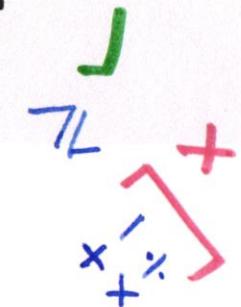


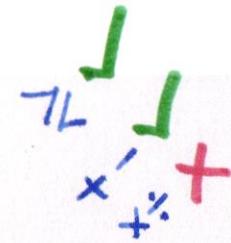
- Nominal data only with two categories.
- **Example:** Winner & Loser.
- **Note:** The **order** is **NOT** important such as Male/Female or Female/Male

## Ordinal Data



- Data values have meaningful order or ranking.
- Magnitude between successive values is not known.
- **Example:** Economic Status

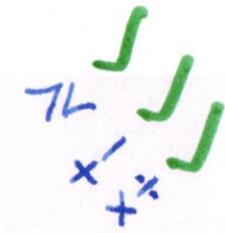
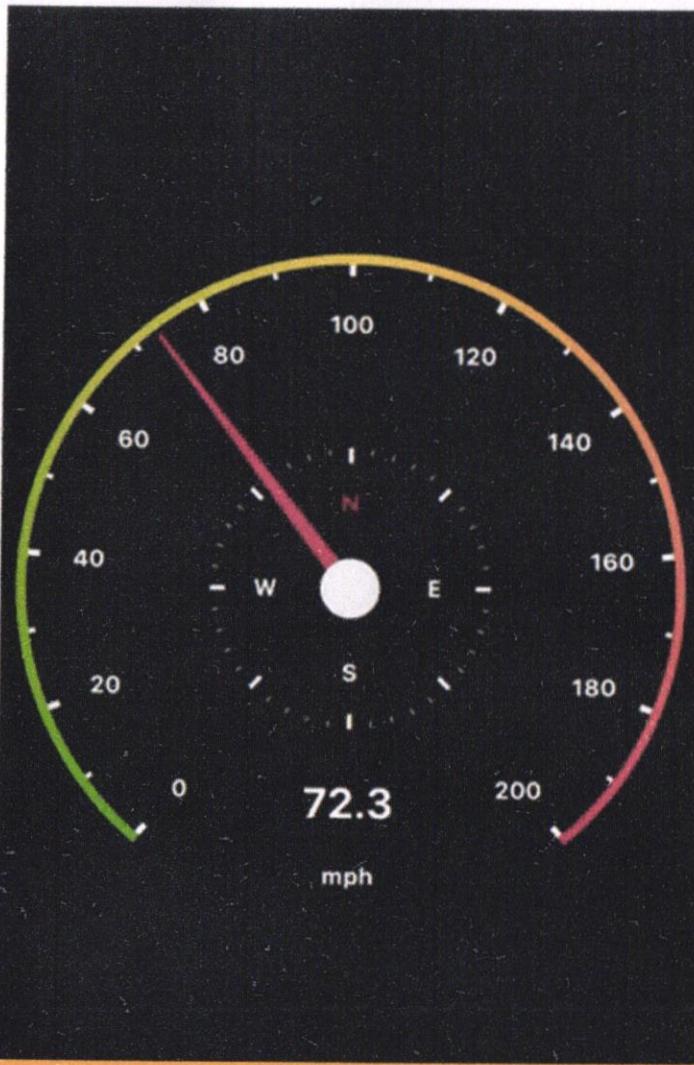




## Numeric Data – Interval Scaled



- It is quantitative in nature with non-inherent zero-point.
- Can be **discrete** or **continuous** in nature.
- Measured on scaled of equal units.
- Can be Positive, Zero or Negative.
- Allows quantify the difference between values.
- **Example:** Temperature



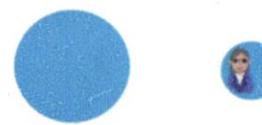
## Numeric Data – Ratio Scaled

- It is quantitative in nature with true zero-point.
- Can be **discrete** or **continuous** in nature.
- Measured on scaled of equal units.
- Can be Positive, Zero or Negative.
- Allows quantify the difference between values.
- **Example:** Speedometer

# Data Quality Issues

|         | A | B  | C | D  | E  | F  | G  | ...  | Attributes. |
|---------|---|----|---|----|----|----|----|------|-------------|
| 1       | N | id | L | Gm | 10 | 12 | .. | ...  |             |
| 2       | P |    | □ | □  | □  | □  |    | .... |             |
| 3       |   |    |   |    |    |    |    |      |             |
| 4       |   |    |   |    |    |    |    |      |             |
| 5       |   |    |   |    |    |    |    |      |             |
| :       |   |    |   |    |    |    |    |      |             |
| entries |   |    |   |    |    |    |    |      |             |

- Duplicity ←
- Inconsistency ←
- Correctness ←
- Timeliness ←
- Missing Values ←



*Exit Polls.*

## Population vs Sample

Population



*Random.* ↓ *Representative*  
Sample



- **Population:** Entire group that you want to draw conclusions about.
- **Sample:** Specific group that you will collect data from.
- The size of the **sample** is always less than the total size of the **population**

## Example: Population vs Sample



- Company A wants to **buy 50,000 batteries.**
- Only buys if no more than 1% are defective.
- Testing all batteries is **expensive.**
- Instead select few samples (500) and test them.
- The results can show estimate of the percentage of defective batteries in population.

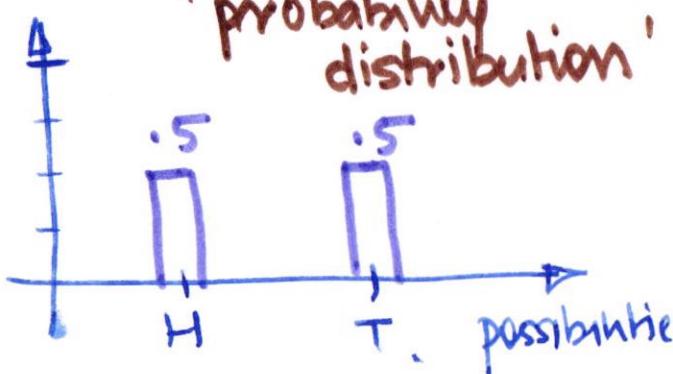
Those variables

whose value is unknown  
until the experiment is  
done.

but

we know something  
about it.

probabilities



Random variable

Probability

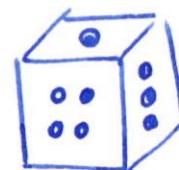
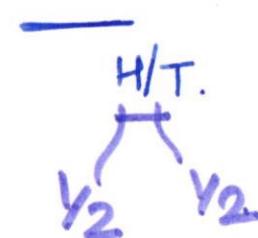
chances of  
occurring

unbiased

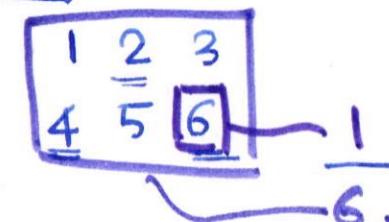
(H)

$$P(H) = \frac{1}{2}$$

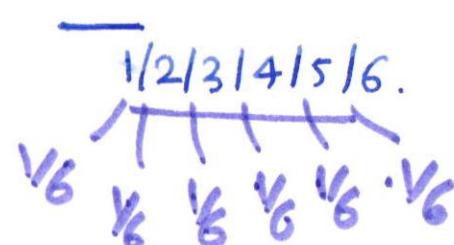
# fav. outcomes  
Total # outcomes



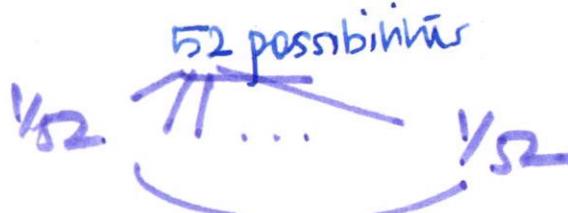
$$P(6) = \frac{1}{6}$$



$$P(\text{even}) = \frac{3}{6} = \frac{1}{2}$$

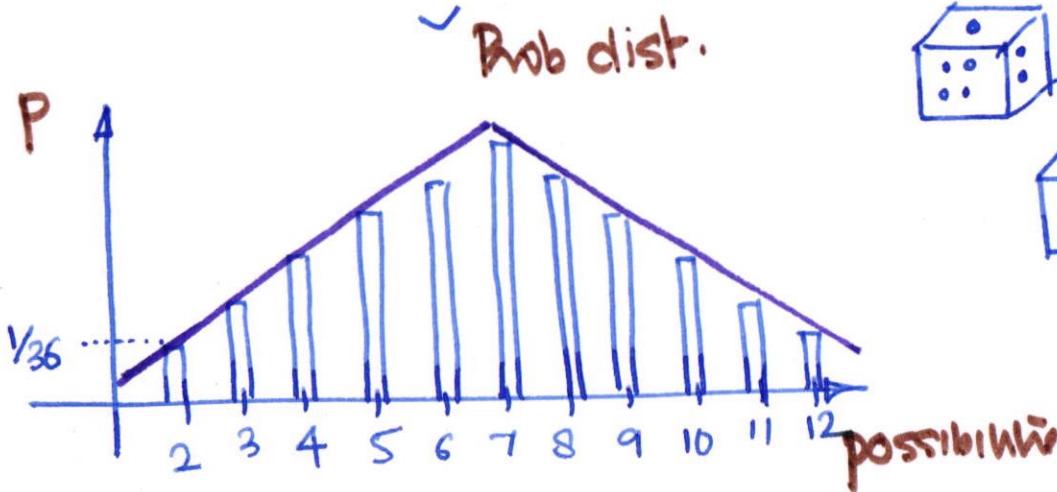


cards.

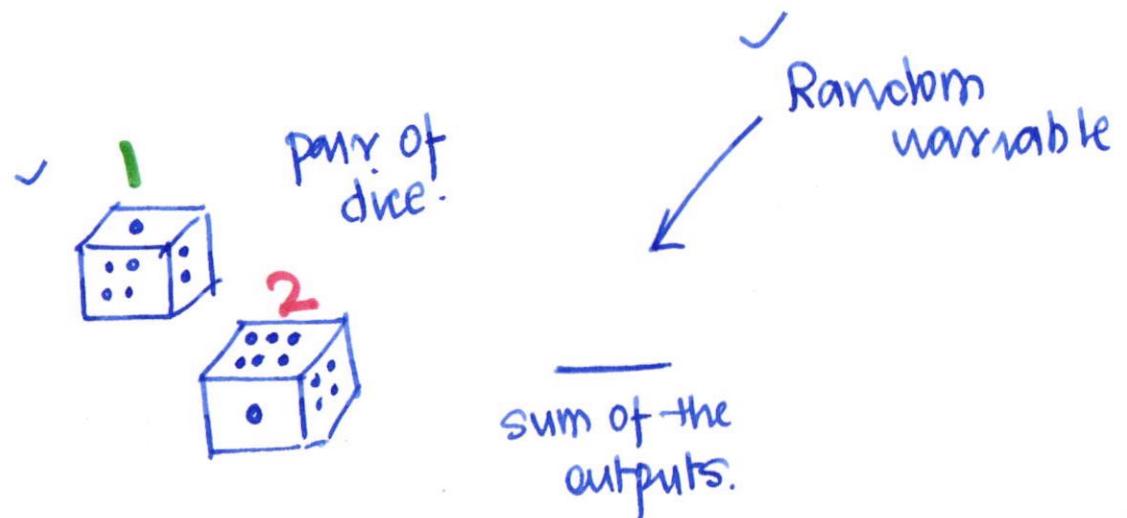


$$P(H) = \frac{13}{52} = \frac{1}{4}$$

$$P(\text{king}) = \frac{4}{52} = \frac{1}{13}$$



✓ prob. distributions  
can be different

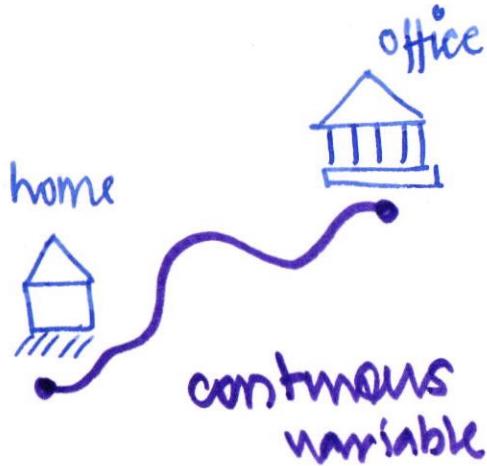


[2, 3, 4, 5, 6, 7, 8, 9,  
10, 11, 12]

outcome 1

outcome 2

|   |   |   |   |   |   |    |
|---|---|---|---|---|---|----|
|   | 1 | 2 | 3 | 4 | 5 | 6  |
| 1 | 2 | 3 | 4 | 5 | 6 | 7  |
| 2 | 3 | 4 | 5 | 6 | 7 | 8  |
| 3 | 4 | 5 | 6 | 7 | 8 |    |
| 4 | 5 | 6 | 7 | 8 |   |    |
| 5 | 6 | 7 | 8 |   |   |    |
| 6 | 7 | 8 |   |   |   | 12 |

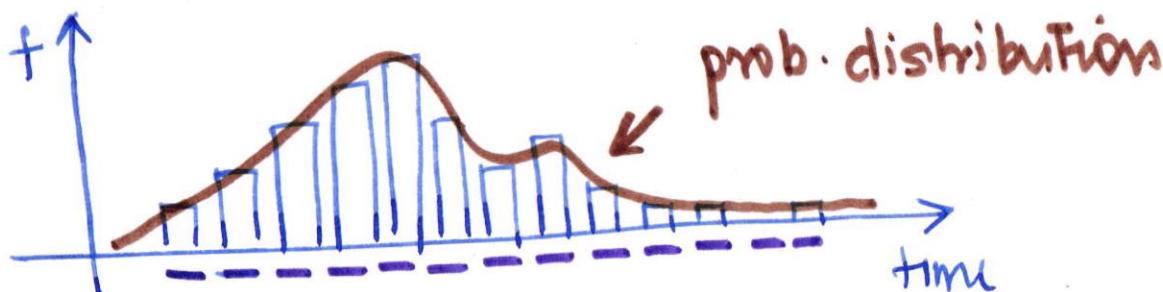
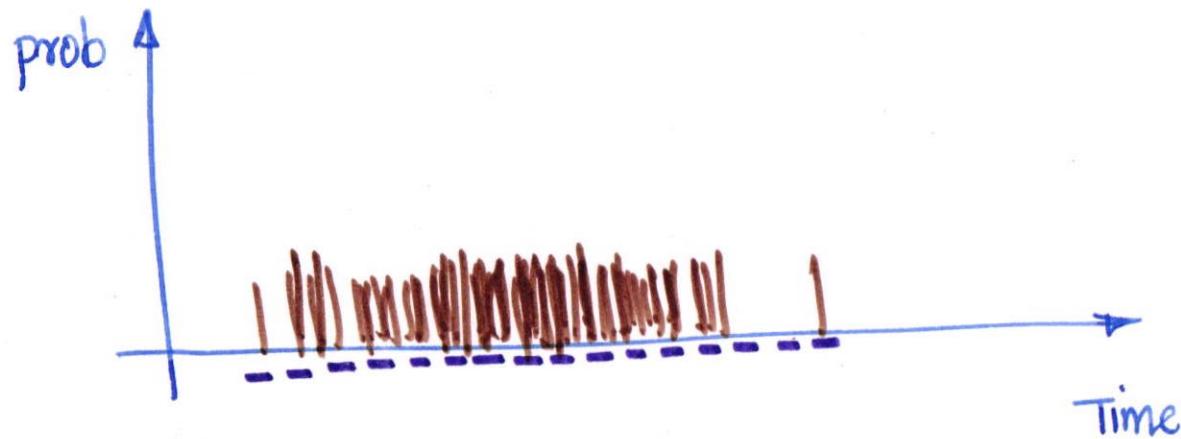


$$\text{distance}_{\text{Home-office}} = \underline{\quad} \text{RV} \times \underline{\quad}$$

$$\text{Time} = \underline{\quad}$$

$$\begin{matrix} \text{Time:} \\ \underline{20.02} & \underline{22.56} \\ \underline{21.31} & \underline{21.37} \\ \underline{51.01} & \underline{15.34} \end{matrix}$$

traffic,  
speed, mode of  
trans,  
weather, 'barrat'



$$\begin{matrix} [20:00 - 20:15, \\ 20:16 - 20:30] \end{matrix}$$



Exit Polls.

130 voters.

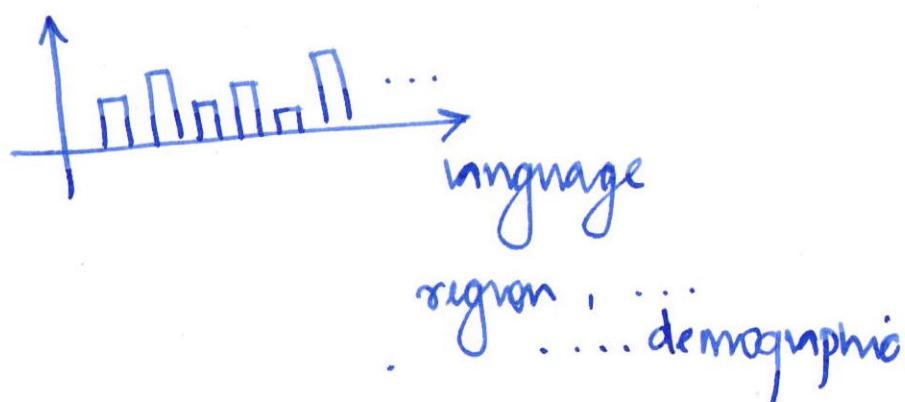


Sample

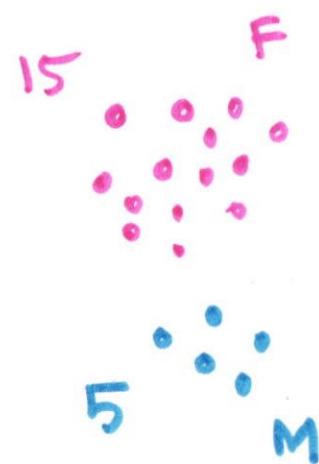


Representative

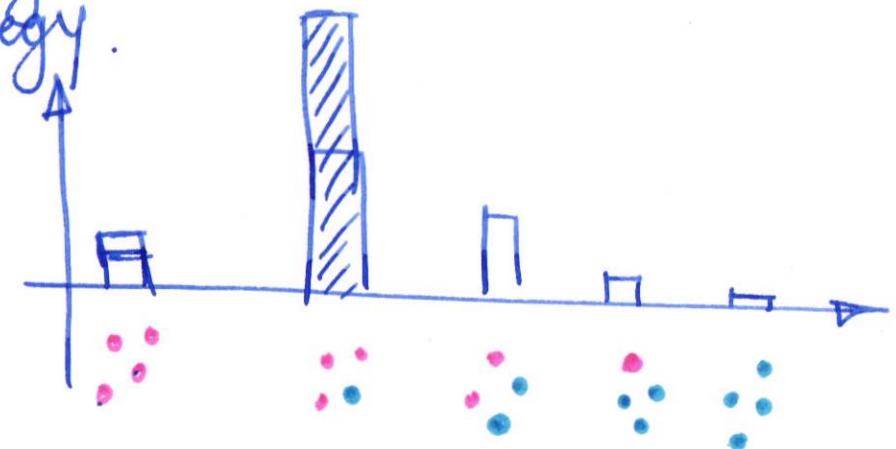
$\text{prob. dist. of SAMPLE} \approx \text{prob. dist. of POPULATION}$



region ... demographic

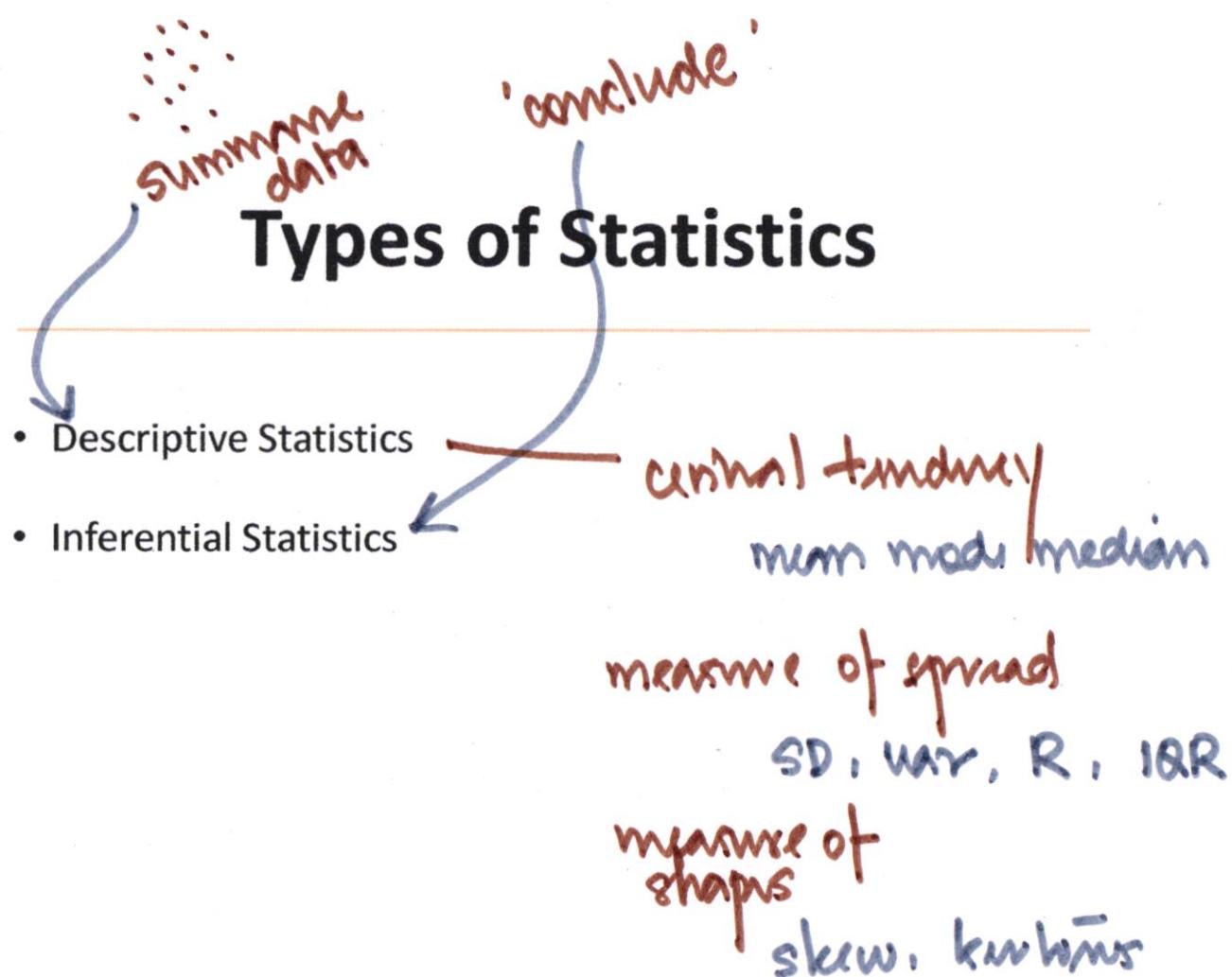
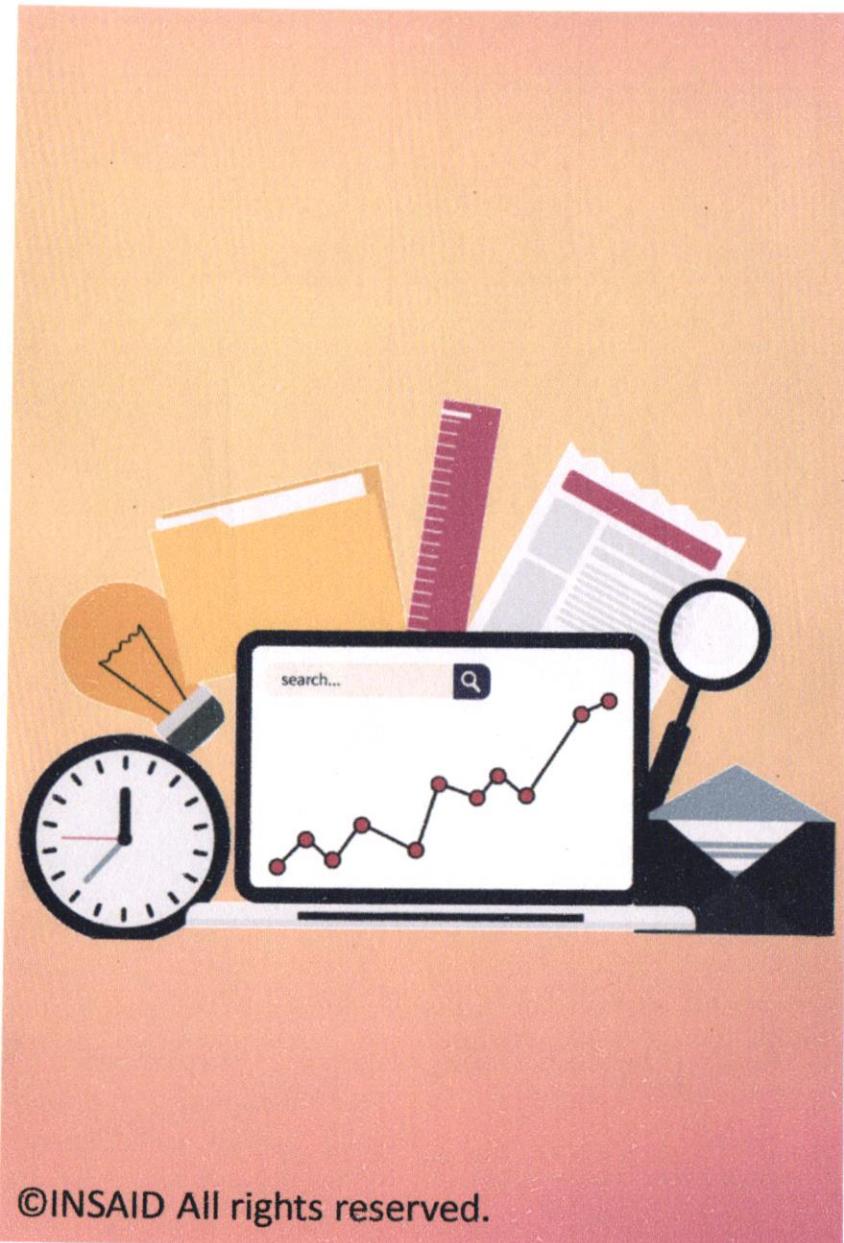


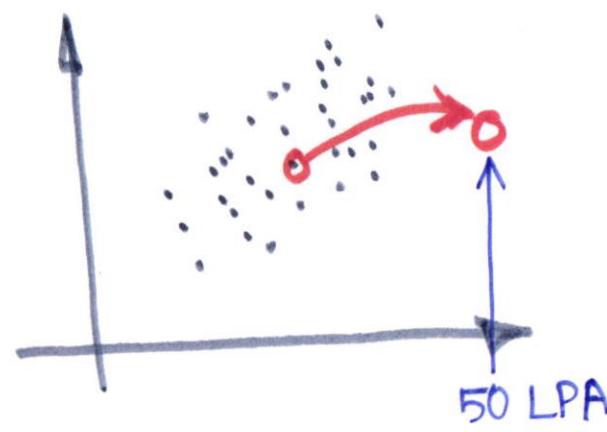
sampling strategy.



4 people  
at RANDOM.

Randomly with each person  
being chosen  
equally likely .



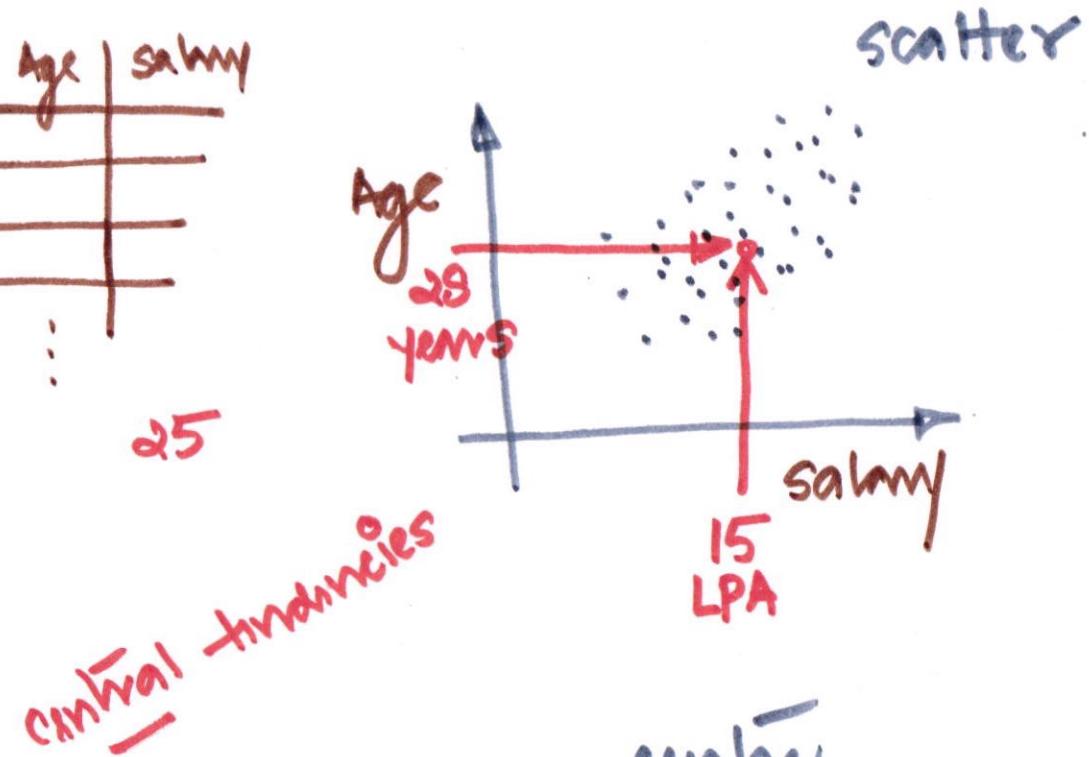


$$\sum_{i=1}^n \text{Variable}_i$$

problem -  
if there is an outlier, avg is not representative

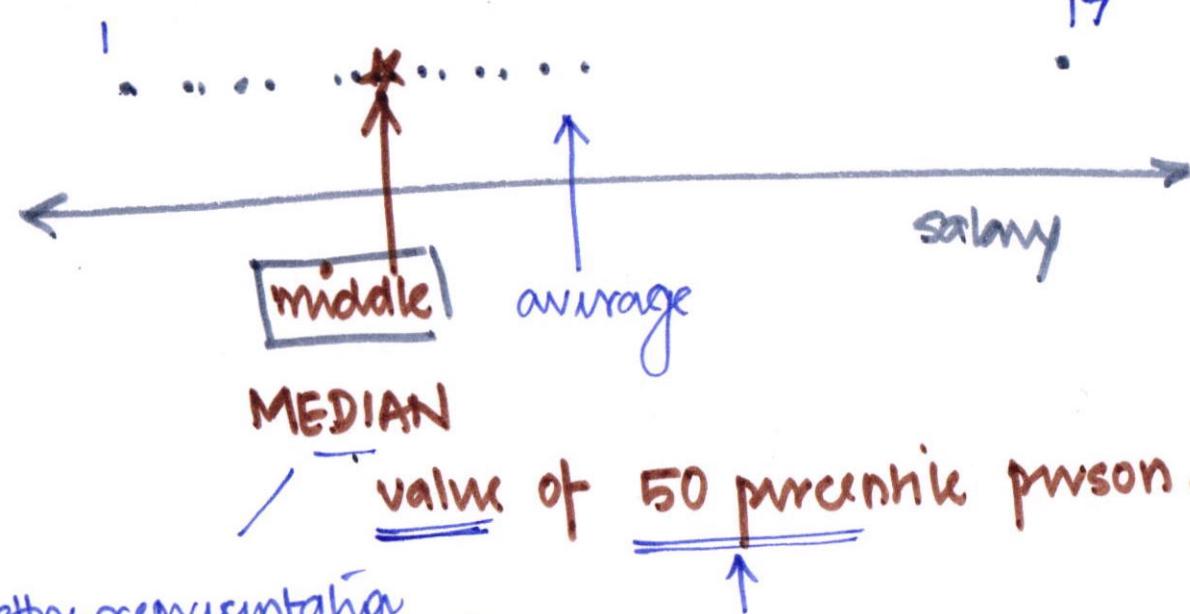
...:

| Age | salary |
|-----|--------|
| 20  | 10 LPA |
| 21  | 12 LPA |
| 22  | 15 LPA |
| 23  | 18 LPA |
| 24  | 20 LPA |
| 25  | 22 LPA |
| 26  | 25 LPA |
| 27  | 28 LPA |
| 28  | 30 LPA |
| 29  | 35 LPA |
| 30  | 40 LPA |
| 31  | 45 LPA |
| 32  | 50 LPA |



Average  $\leftrightarrow$  mean

$$\frac{\sum_{i=1}^n \text{salary}_i}{n}$$



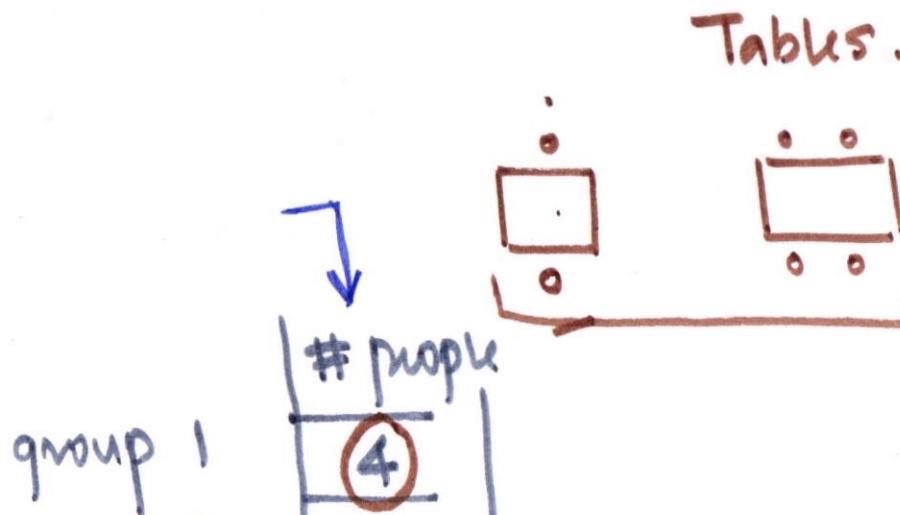
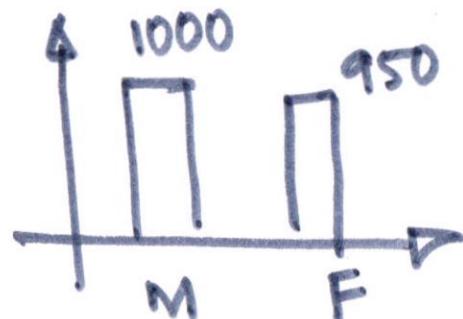
better representation  
than average, if there  
is an outlier.



mode

## MODE

categorical variables



group

| # people |
|----------|
| 4        |
| 3        |
| 1        |
| 6        |
| 7        |
| 4        |
| 3        |

most frequent

4.25

## MODE

value which is most frequent.

$$\text{mode}(\text{Grw}) = M$$

# Descriptive Statistics

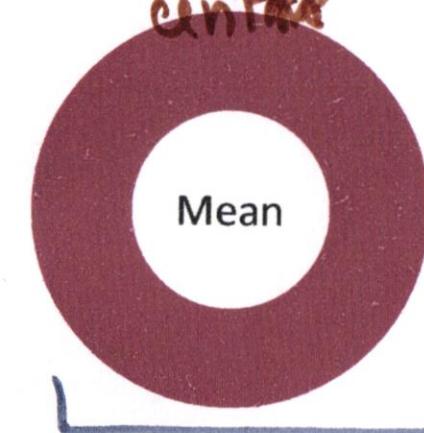
- Measures of Central Tendency
- Measures of Spread
- Measures of Shape



# Measures of Central Tendency

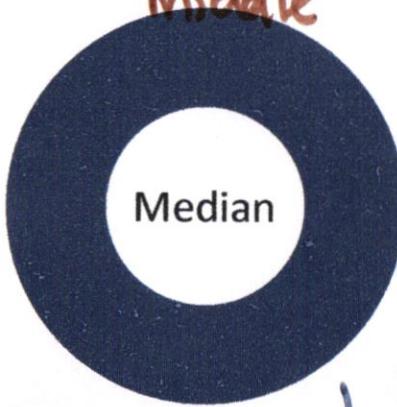
$$\frac{\sum_{i=1}^n x_i}{n}$$

center



value at  
50<sup>th</sup> percentile

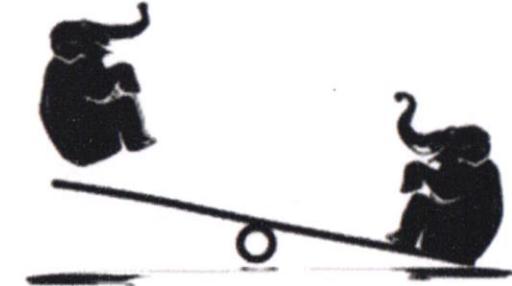
Middle



most frequent

most frequent

Mode





Rs. 5 Lakhs



Rs. 6 Lakhs



Rs. 7 Lakhs



Rs. 8 Lakhs

## Mean

- It is the **average** of a set of data.
- **Example:** Average salary of employees
- $\text{Mean} = (5 + 6 + 7 + 8) \div 4 = 6.5 \text{ Lakhs}$



Rs. 5 Lakhs



Rs. 6 Lakhs



Rs. 7 Lakhs



Rs. 25 Lakhs

## Mean (Outliers Effect)

- Mean provides a misleading balance point because of an outlier.
- $\text{Mean} = (5 + 6 + 7 + 25) \div 4 = \underline{\underline{10.7}} \text{ Lakhs}$



30 Kg



35 Kg



38 Kg



40 Kg



42Kg

## Median

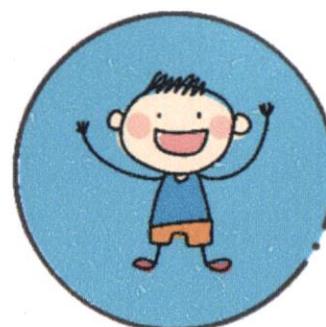
- It is the **middle value** in a set of ordered data values.
- **Example:** Weight of kids.



30 Kg



35 Kg



38 Kg



40 Kg



75 Kg

## Median (Outliers Effect)

- Median is a positional average and refers to the middle value in a distribution.
- **Example:** Weight of kids.



Female



Female



Male



Male



Male

Mode = 2

Mode = 3

## Mode

- It is the most frequently occurring value in the dataset.
- **Example:** In above diagram, Male is occurring the most.



## Pizza Table Example

---

- Deciding the seating arrangement of a restaurant by figuring out the most frequently occurring group size.
- A sample of 20 groups is selected at random:
- People = {2, 4, 1, 2, 3, 2, 4, 2, 3, 6, 8, 4, 2, 1, 7, 4, 2, 4, 4, 3}
- There are 2 modes, each occurring six times – 2, 4
- Basis the above, the manager will decide on 2 seater and 4 seater tables being kept in the patio.

hope - the center represents the group.

## Central Tendency Use Cases

- Understand data spread around the center in Data Analysis
- Replace missing values in Data Analysis
- Showcase performance of players across their career
- Study salary spread across the company
- & Many More...

1

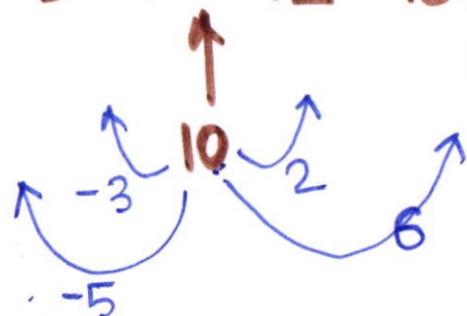


2



~~avg age~~ missing

5 7 12 16.



$$\begin{aligned}\text{average} &= \frac{5+7+12+16}{4} \\ &= 10 \\ &= \mu \text{ (cmu)}\end{aligned}$$

$$\begin{aligned}\text{average mistake} &= \frac{(-3) + (-5) + (2) + (6)}{4} \\ &= 0\end{aligned}$$

~~$$\text{avg. // } \neq$$~~
$$\text{avg error} = \frac{\sum_{i=1}^n e_i}{n} \times$$

we represent  
the data with A B C D E  
**MEAN**

on an average  
how much are we  
wrong, in saying  
centre is representation -ve

$$\begin{aligned}\text{avg. abs. Error} &= \frac{\sum_{i=1}^n |e_i|}{n} \\ &= \frac{3+5+2+6}{4} \\ &= \underline{\underline{4}}\end{aligned}$$

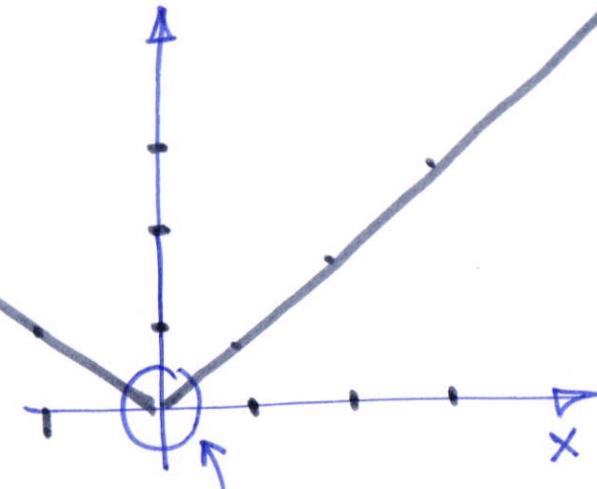
Variance  
 $\sigma^2$

$$\frac{\sum_{i=1}^n (e_i)^2}{n}$$

$$= 18.5 \\ \underline{(\text{Lak Rs})^2}$$

$$e_i = \mu - x_i \\ \text{mean} - x_i$$

$$y = |x|$$



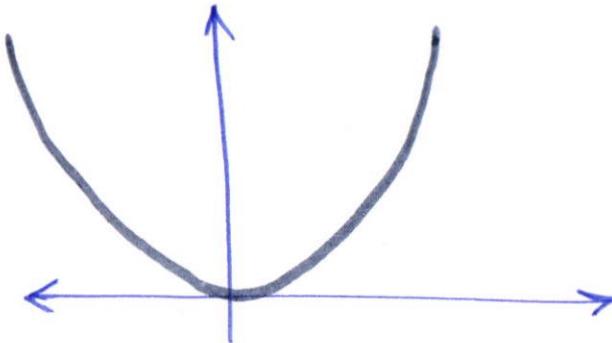
we don't like  
sharp point

lost the  
meaning

standard deviation  
 $\sigma$   
(sigma)

$$\sqrt{\frac{\sum_{i=1}^n (e_i)^2}{n}}$$

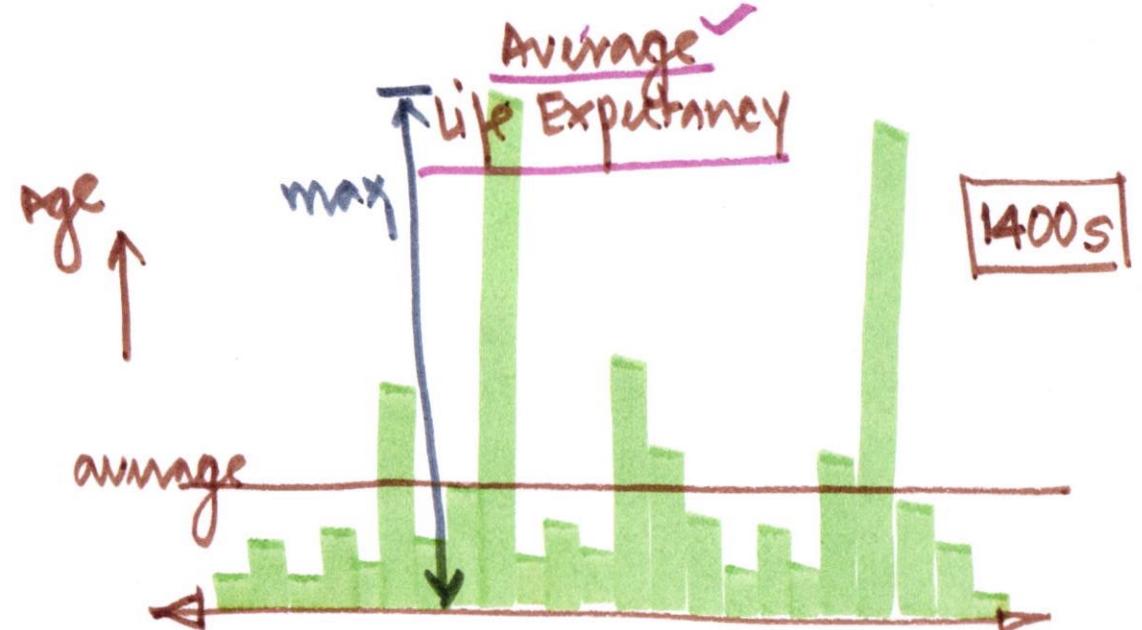
$$= 4.3 \\ \underline{\text{Lakhs}}$$



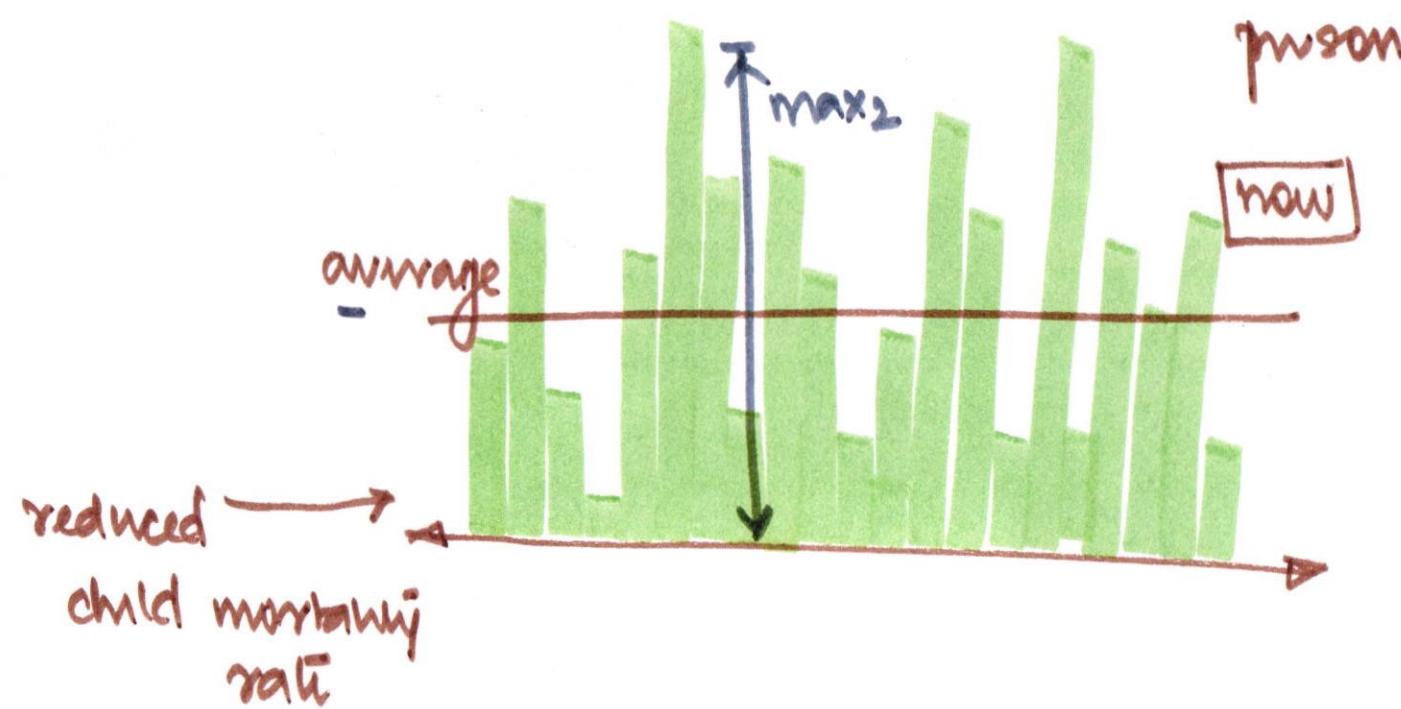
age  
=  
average?  
YES/NO?

$$\max_1 > \max_2$$

AVERAGE

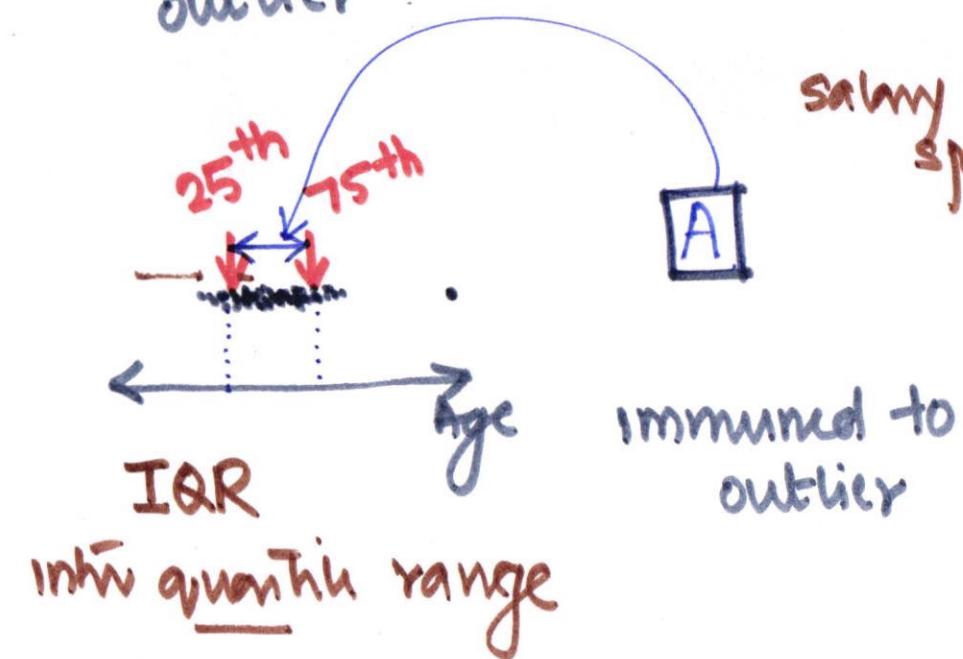
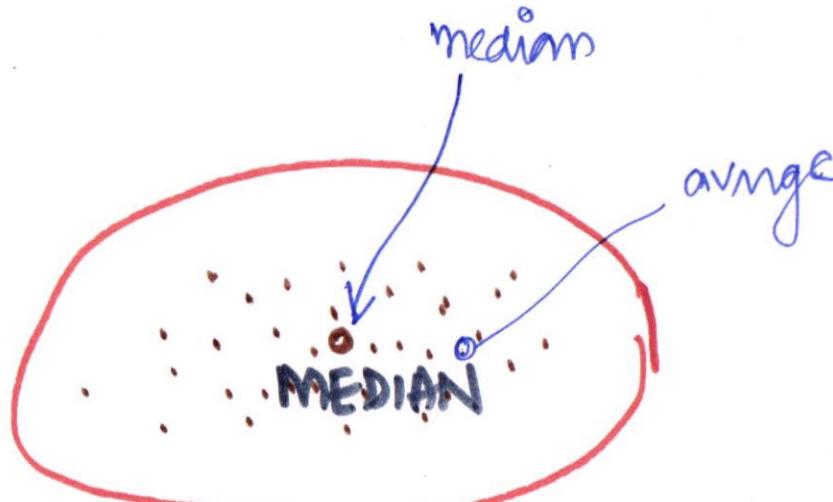
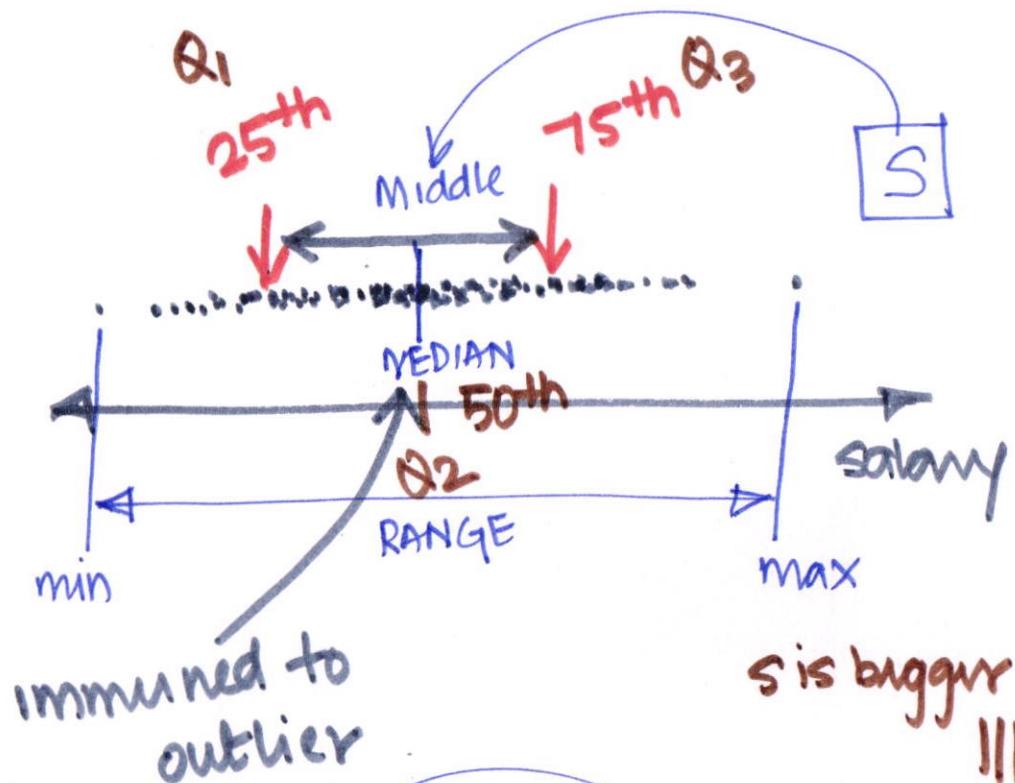


1400s



now

reduced  
child mortality  
rate



Invisible line of demarkation is around the Median.

radius of demarkation  $\# \text{IQR}$

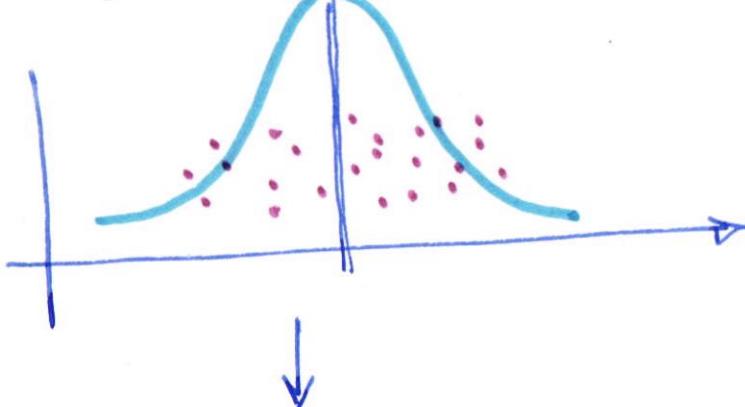
If point under study = far enough from MIDDLE Then we can call it outlier.



\* Average salary = 8.7 LPA

✓ median = 6.9 LPA

IQR =



no outlier



6 LPA

C1

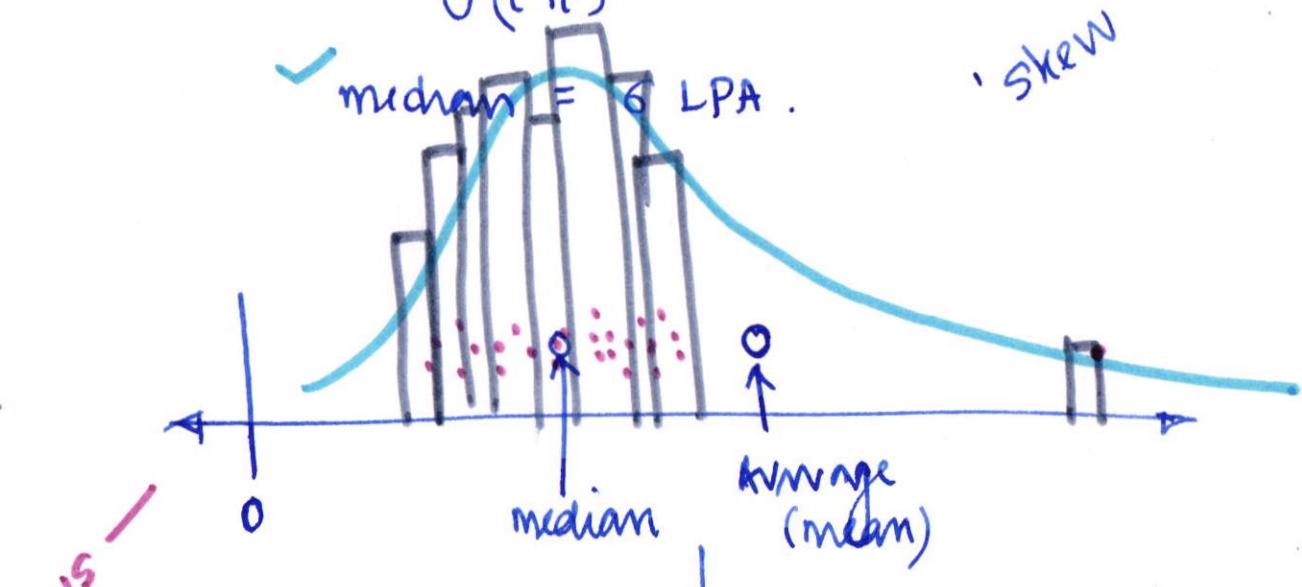


\* Average salary = 10. LPA

✓ median = 6 LPA

(CTC)

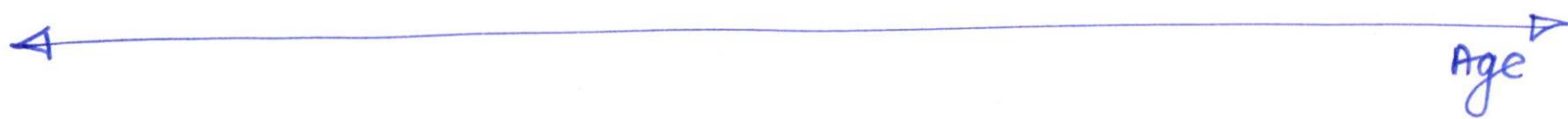
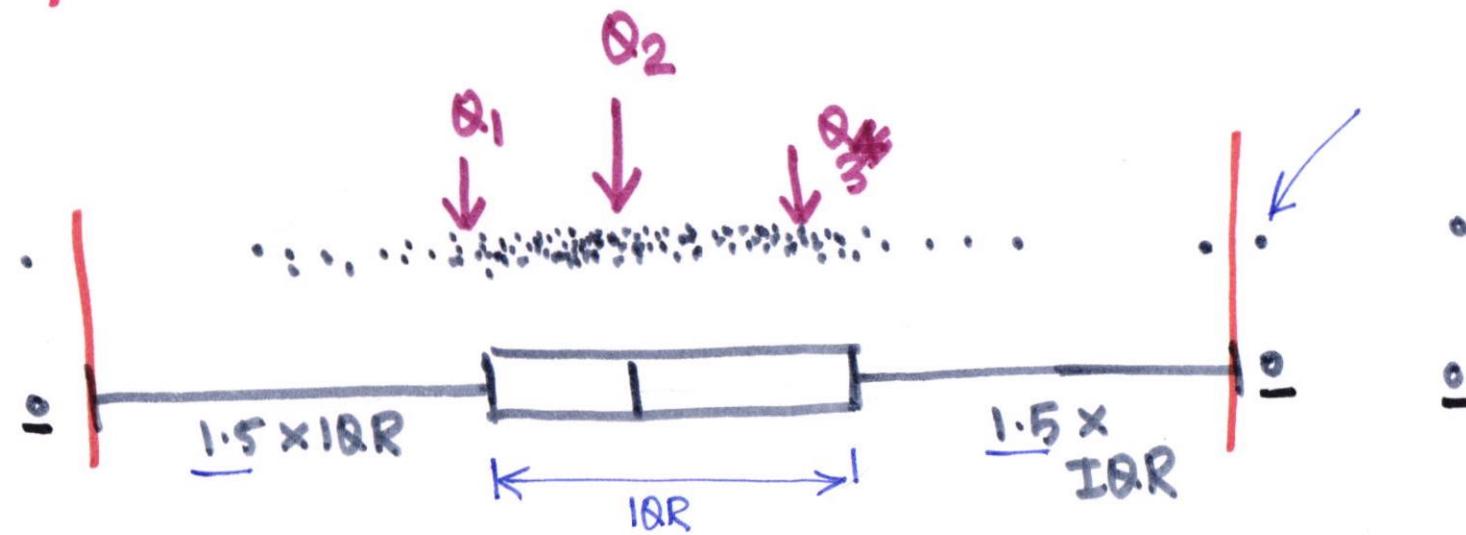
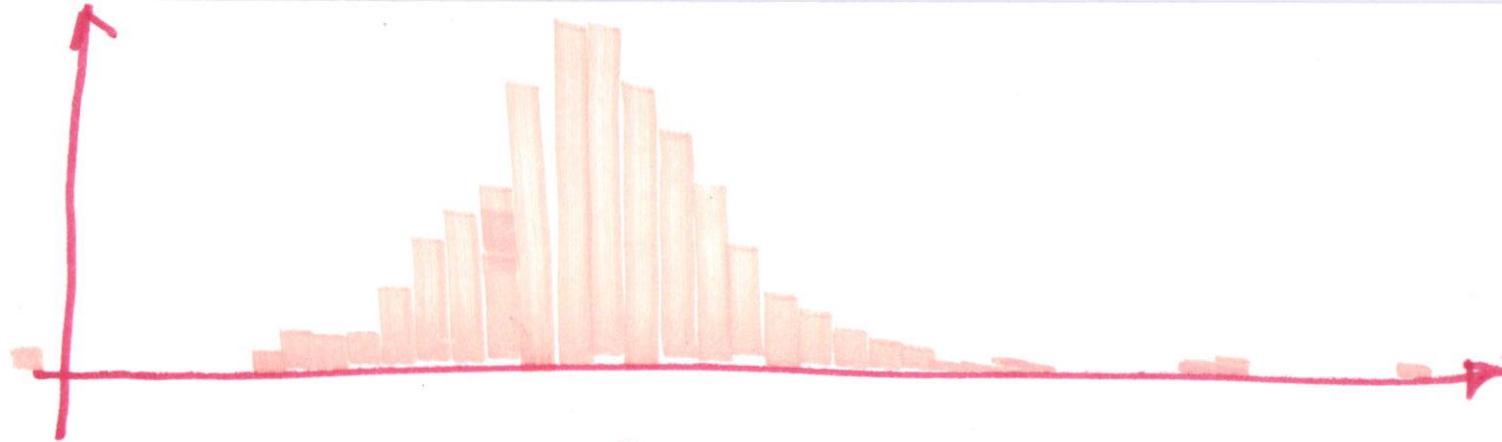
'skew'



average (mean)

There is outlier



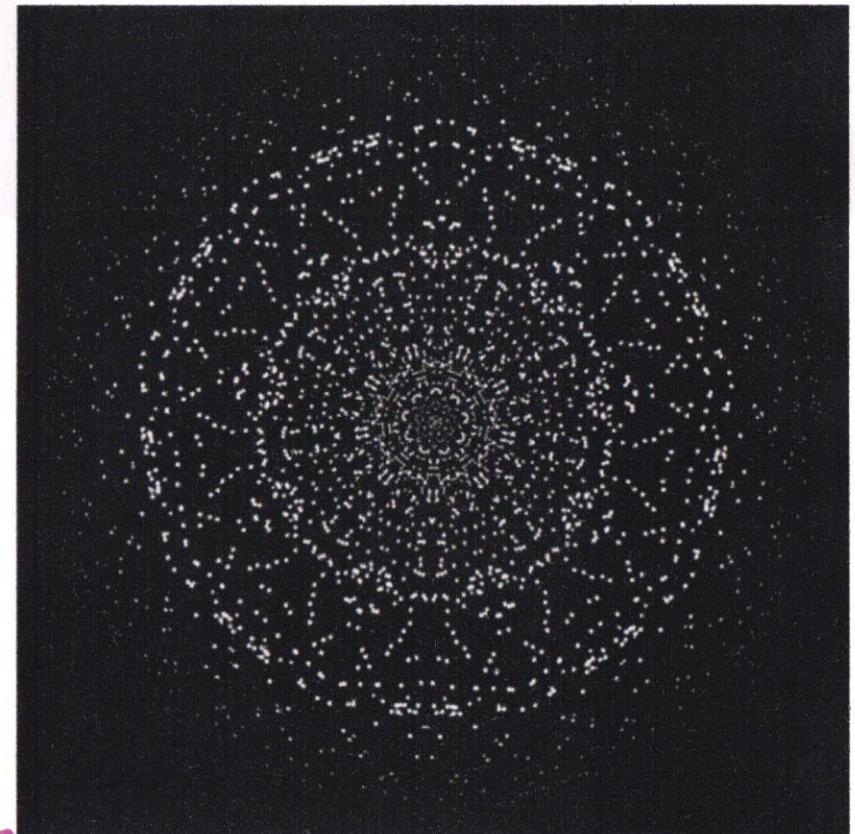


# Measures of Spread

- Range
- Quantiles & Quartiles
- Interquartile Range
- Standard Deviation
- Variance

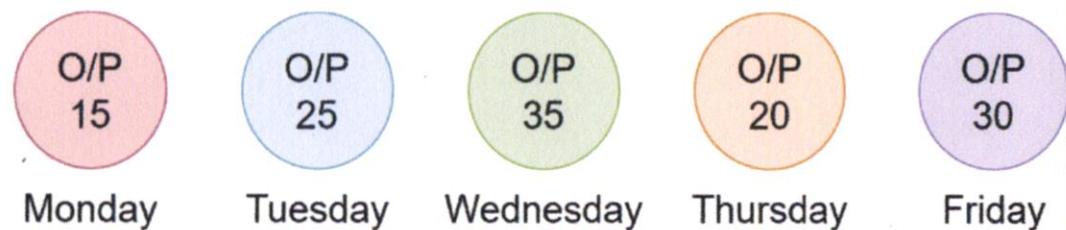
## Positional

- Relm  
of  
**Median**
- = feeling,  
that we  
might have  
outliers
- centre  
mean  
really  
but, affected by outliers



# Range

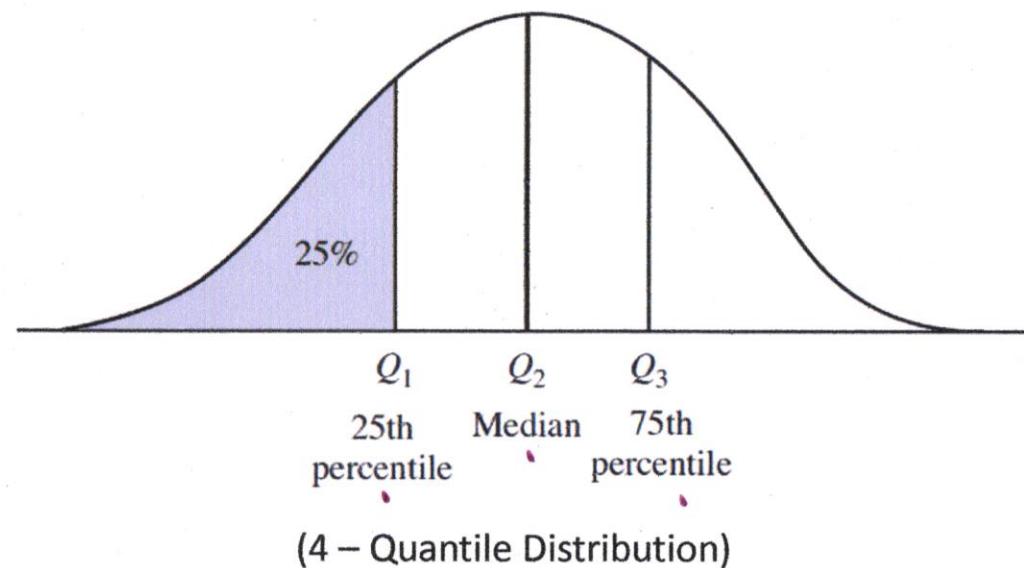
- Difference between the largest and smallest value.
- **Example:** Plant manager are asked to record the plant production output during weekdays.



- Range = Max – Min
- Range =  $35 - 15$
- Range = 20

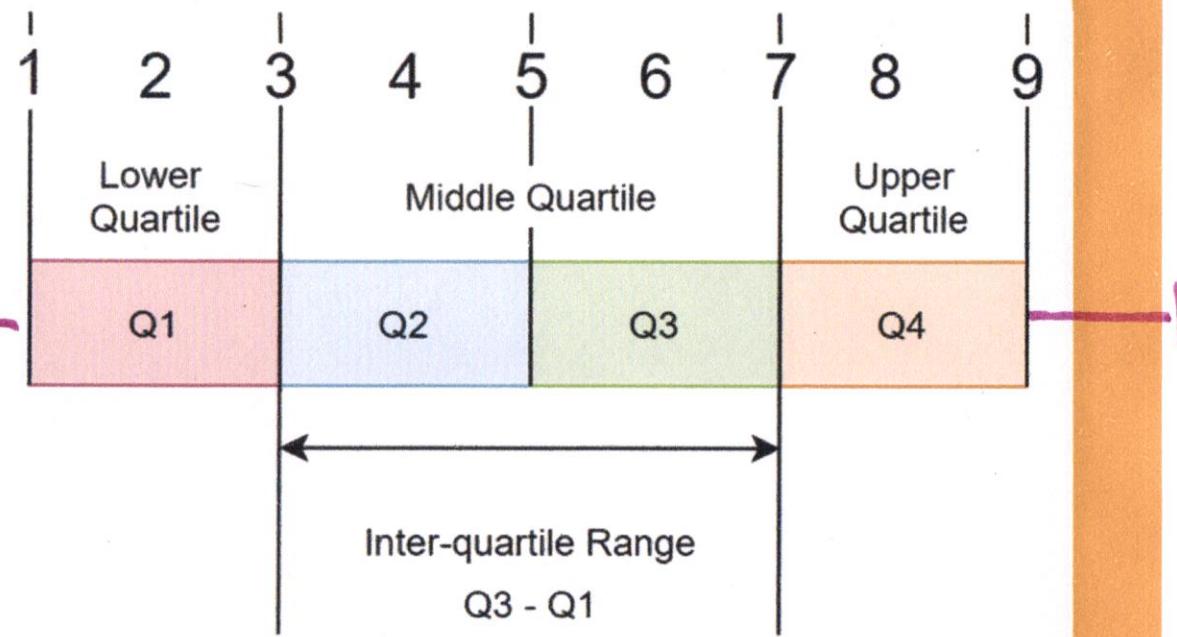
# Quantiles & Quartiles

- **Quantiles:** Data points that split data distribution into equal-size consecutive sets.
- **Quartiles:** Each part split by quantiles.
- **Example:** 4 - Quantile distribution



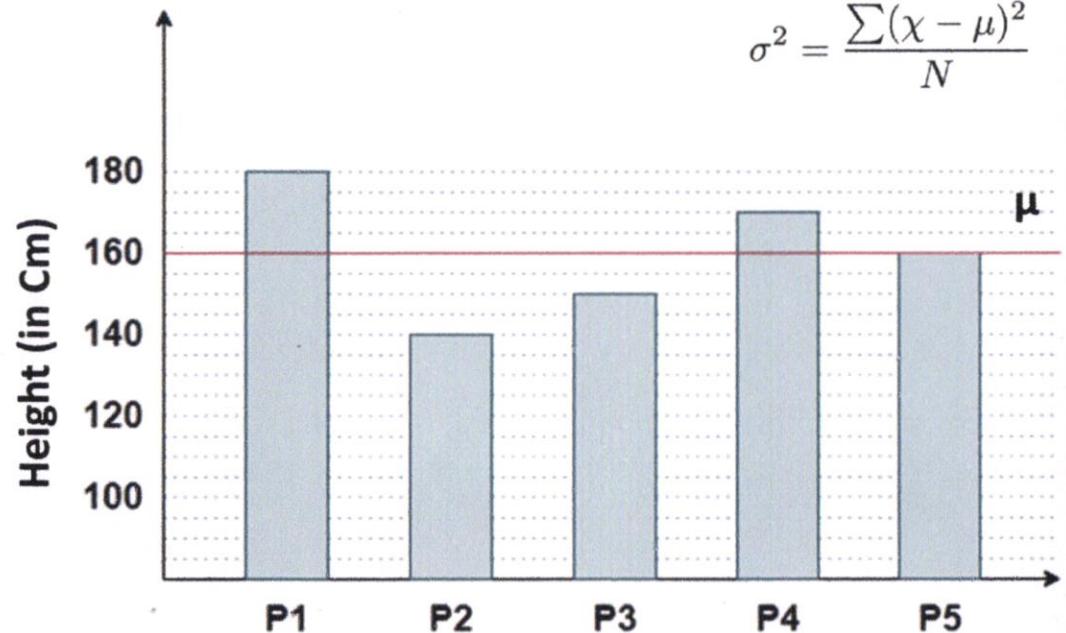
## Interquartile Range

- The distance between the first and third quartiles.
- Gives the range covered by the middle half of the data.



# Variance

- It is the average of the squared distances from the mean.
- **Pros:** Adds weight to outliers, the numbers that are far from the mean.
- **Cons:** It is not easily interpreted.
- **Example:** Height of students.

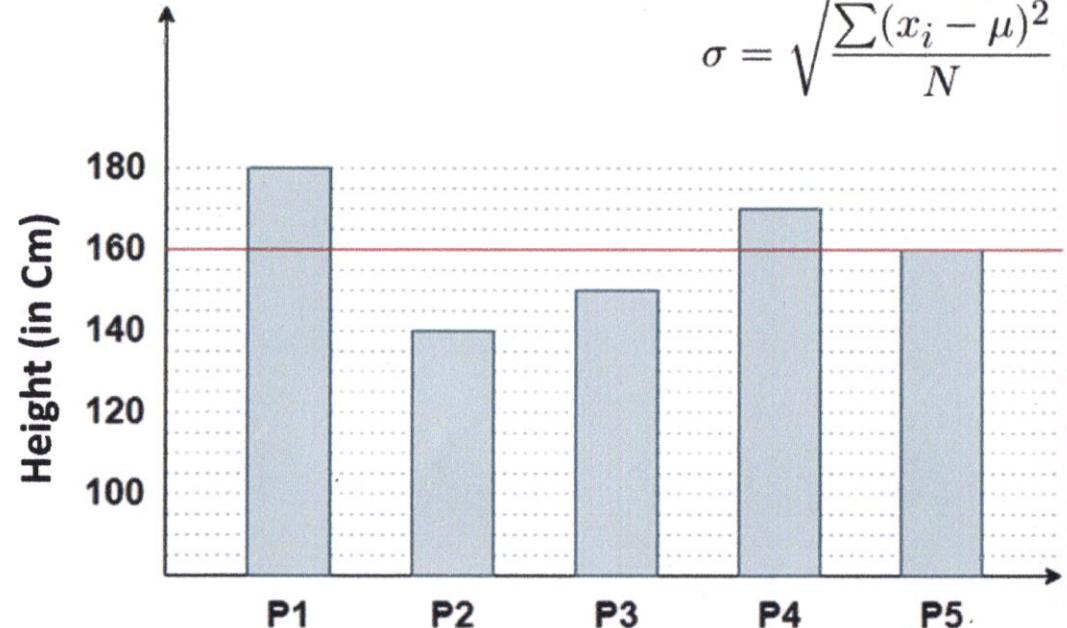


**Solution:**

- $\sigma^2 = (400 + 400 + 100 + 100 + 0) \div 5$
- $\sigma^2 = 200$

# Standard Deviation

- It is the square root of variance.
- Pros: Easy to interpret.
- Cons: Highly affected by the outliers.
- Example: Height of students.



**Solution:**

- $\sigma = \sqrt{[(400 + 400 + 100 + 100 + 0) \div 5]}$
- $\sigma = 14.1$

## *'whining'* *'on and on'* **Measures of Spread Use Cases**

---



- Gives us an idea of how well the mean represents the data.
- It help in figuring out the outliers for better results.
- Standardized the data and interpret with respect to mean.
- Monitor control systems for not normal behavior.
- & Many More...

# Measures of Shape

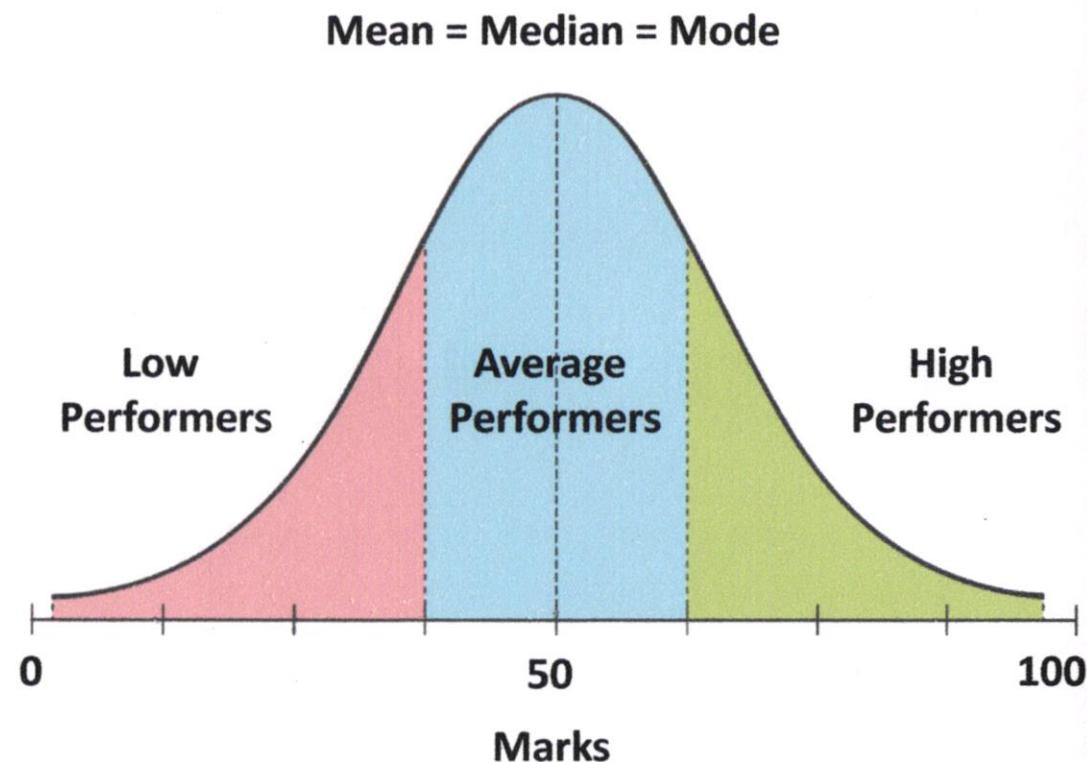
chart it w.r.t Normal distribution

- Symmetry ✓
- Skewness
- Kurtosis



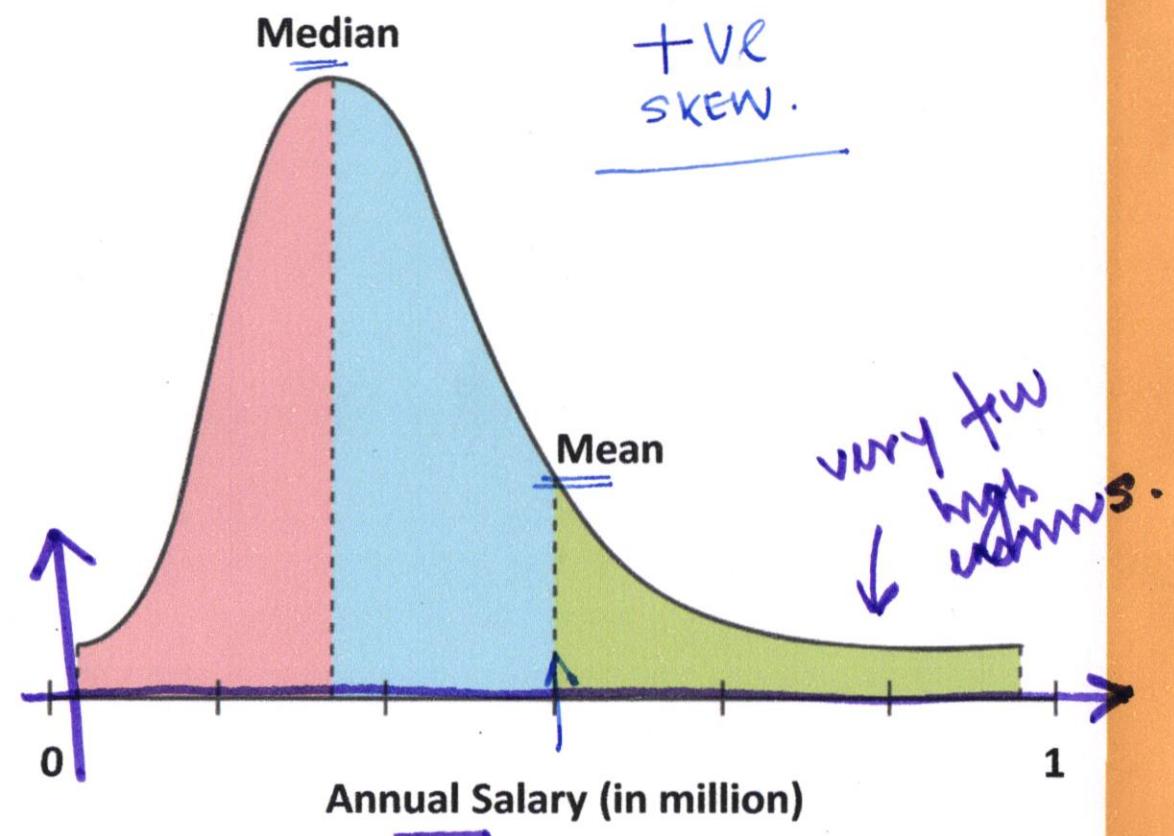
# Symmetry

- Also known as normal distribution.
- The mean, median and mode lie at same point.
- **Example:** Distribution of marks received by 100 students in a math's test.



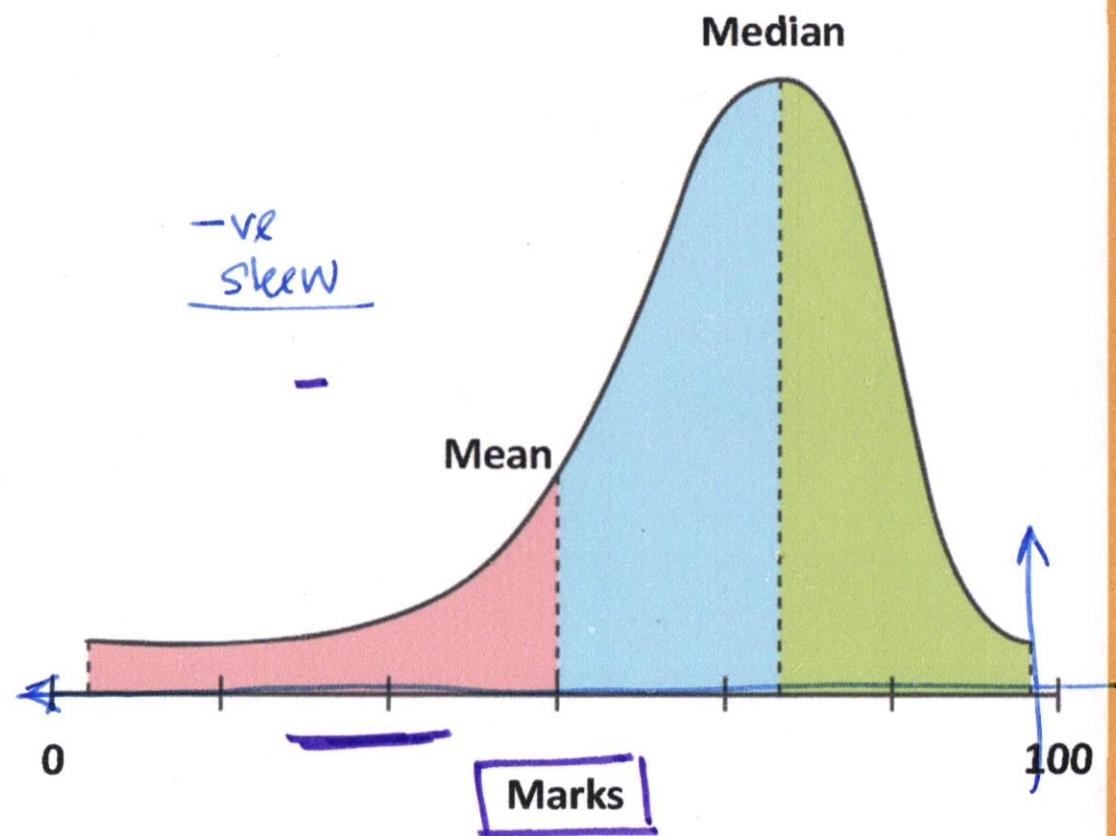
## Skewness - Positive

- Also known as right skewed distribution.
- Mean > Median
- **Example:** Income Distribution in a sample population



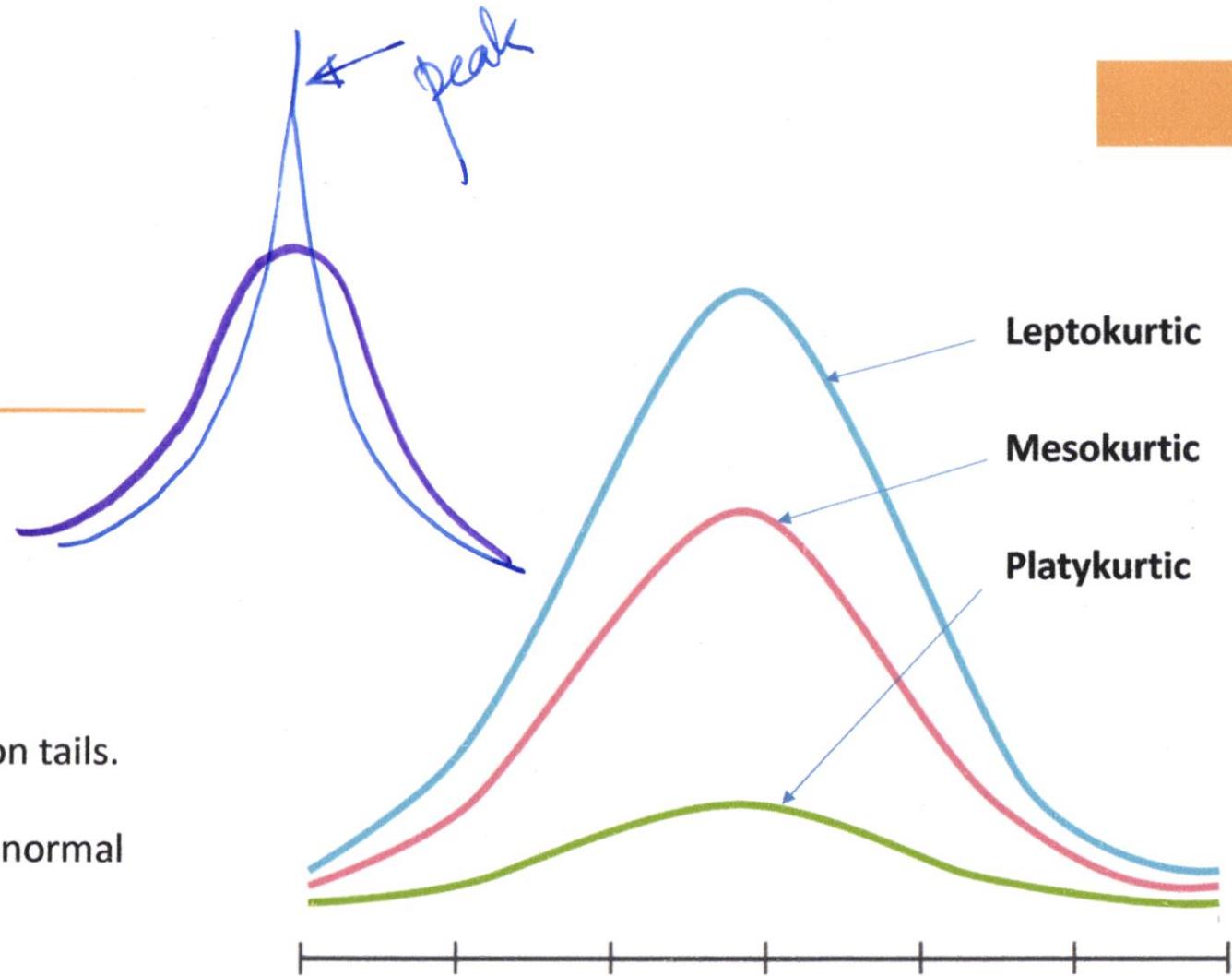
## Skewness - Negative

- Also known as left skewed distribution.
- Mean < Median
- **Example:** Distribution of scores on a very easy test.



## Kurtosis

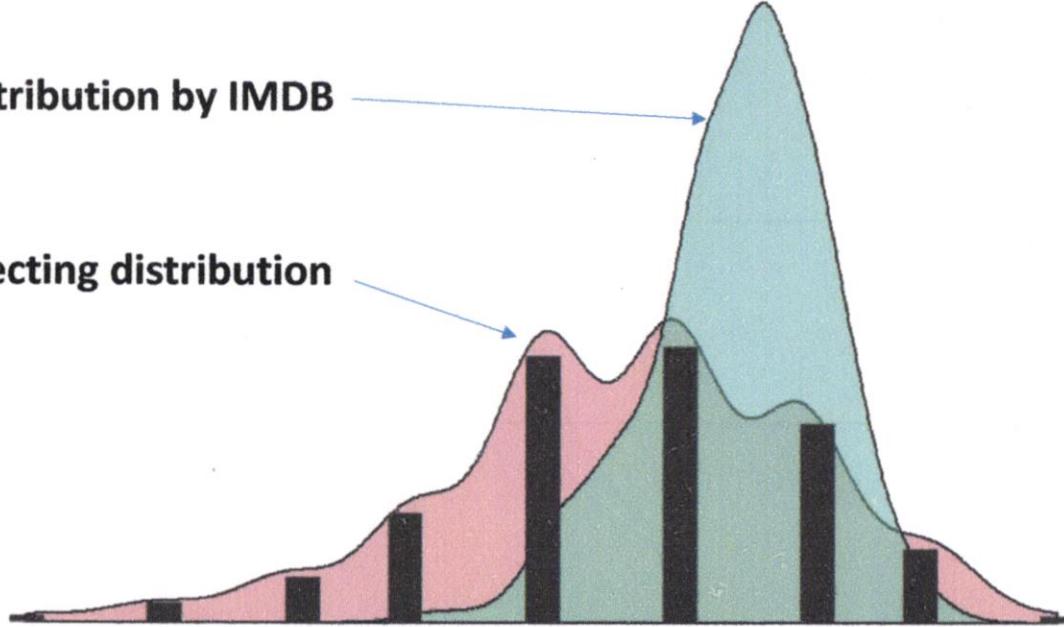
- A statistical measure.
- Defines heaviness of distribution tails.
- Helps in differing the tails with normal distribution.



**Standard distribution by IMDB**

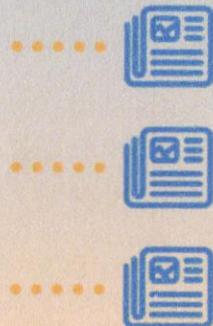
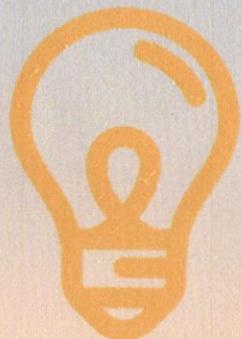
User  
IMDB

**Subjective taste affecting distribution**



## Kurtosis Example

- Comparison of 400+ movie ratings by IMDB and a single user.



Assumption -  
Data is Normal

check , conclusion



- Histograms
- Bar Charts
- Scatter Plots
- Pie Charts
- Box Plots
- & Many More...

## Plots

categorical

axis  
will also  
be  
categorical

nominal

ordinal

interval

ratio

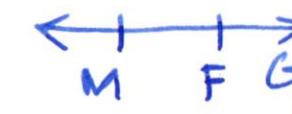
Types of variables



continuous

axis  
will also  
be  
continuous

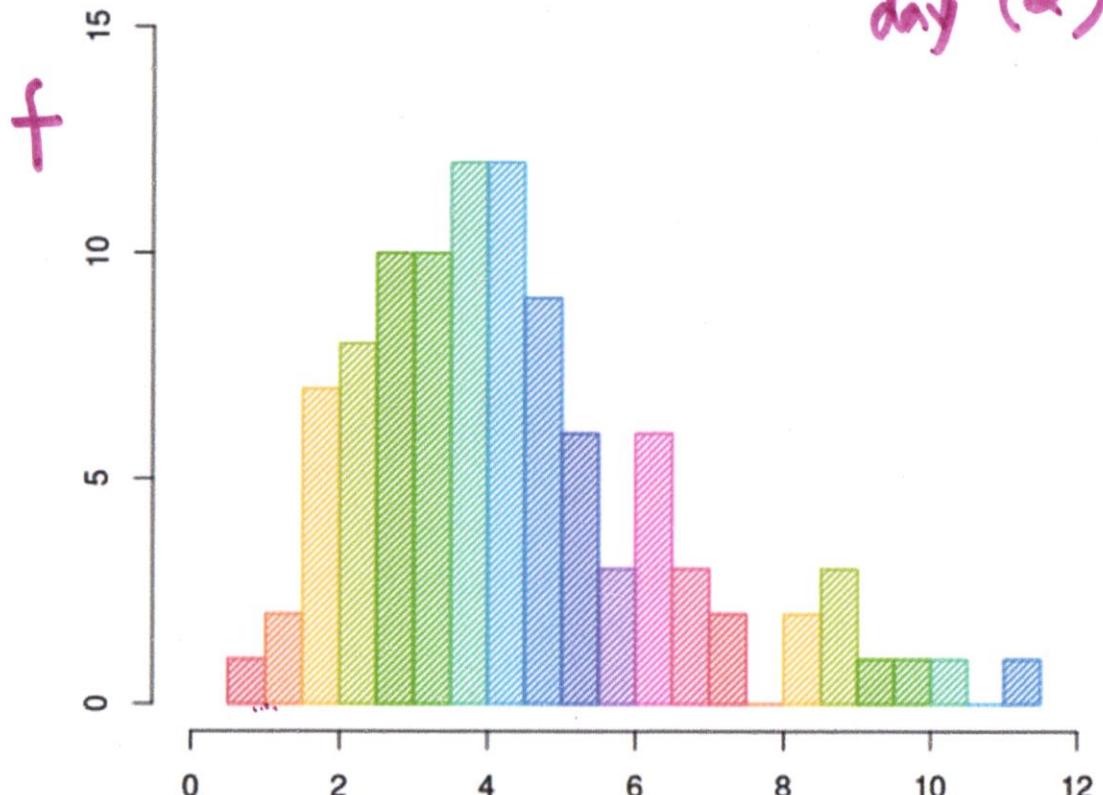
time



# Histogram

*continuous*

- A graphical display of data using bars of different heights.
- Each bar groups numbers into ranges.
- Taller bars show that more data falls in that range.
- Displays the shape and spread of continuous sample data.

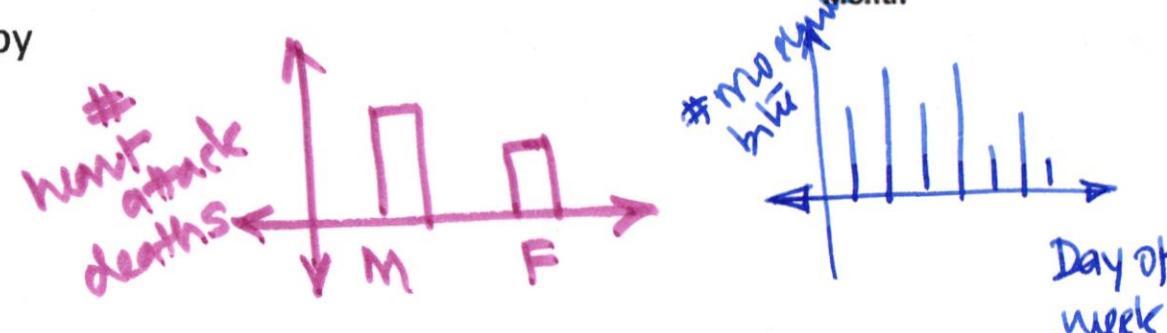
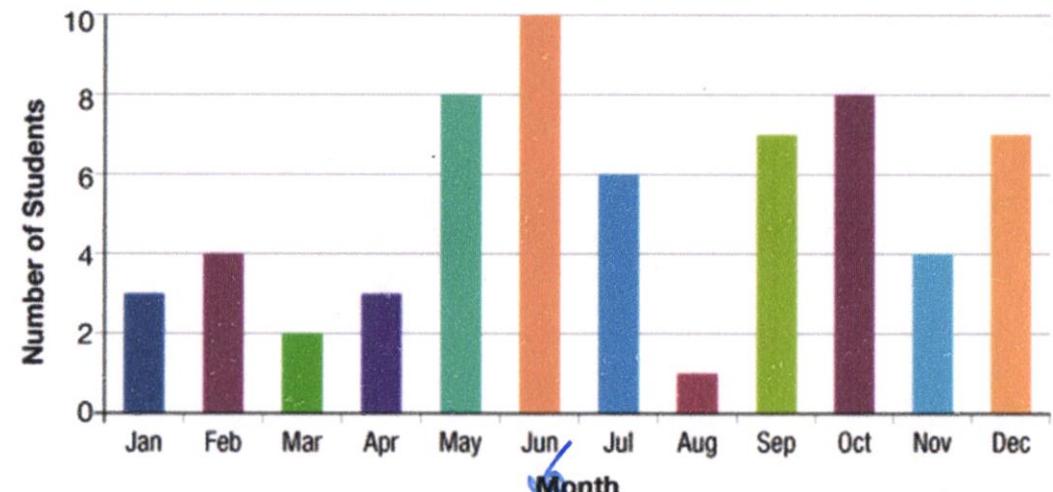


*variable*

# Bar Chart

- Used to show how something changes over time or to compare items.
- Generally, X-axis contains categories and Y-axis contains continuous values.
- Example:** Birthday of Students by Month

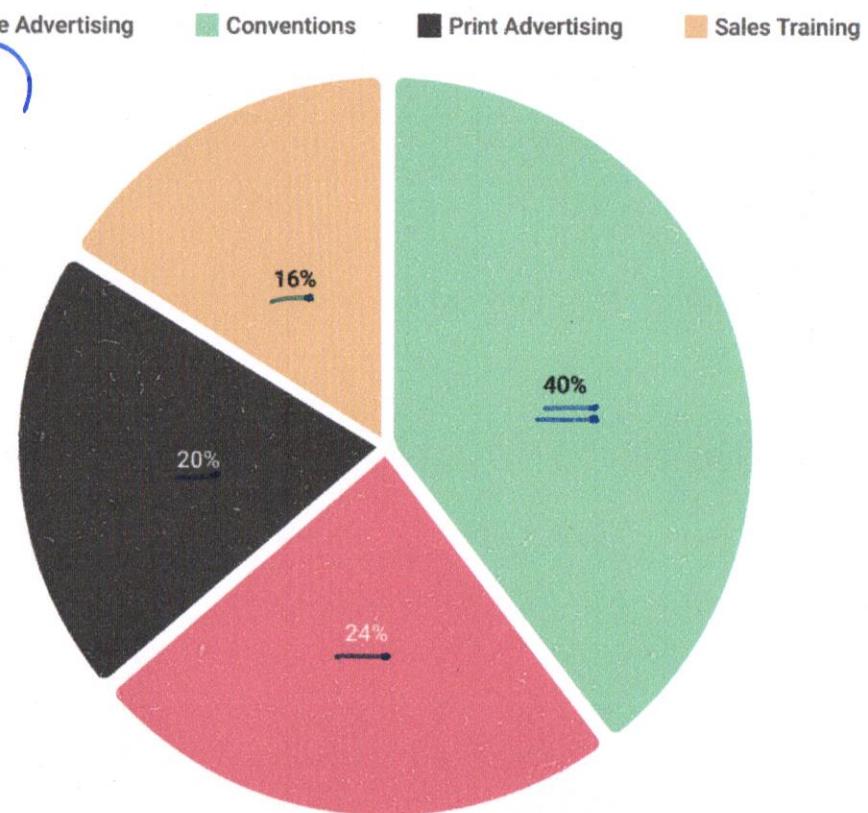
Birthday of Students by Month



# Pie Chart

- A circular statistical graph.
- It is divided into slices to illustrate numerical proportion.
- The arc length of each slice represents proportion to the quantity.
- **Example:** Marketing Budget

*category data*  
*Frequency of category (%)*

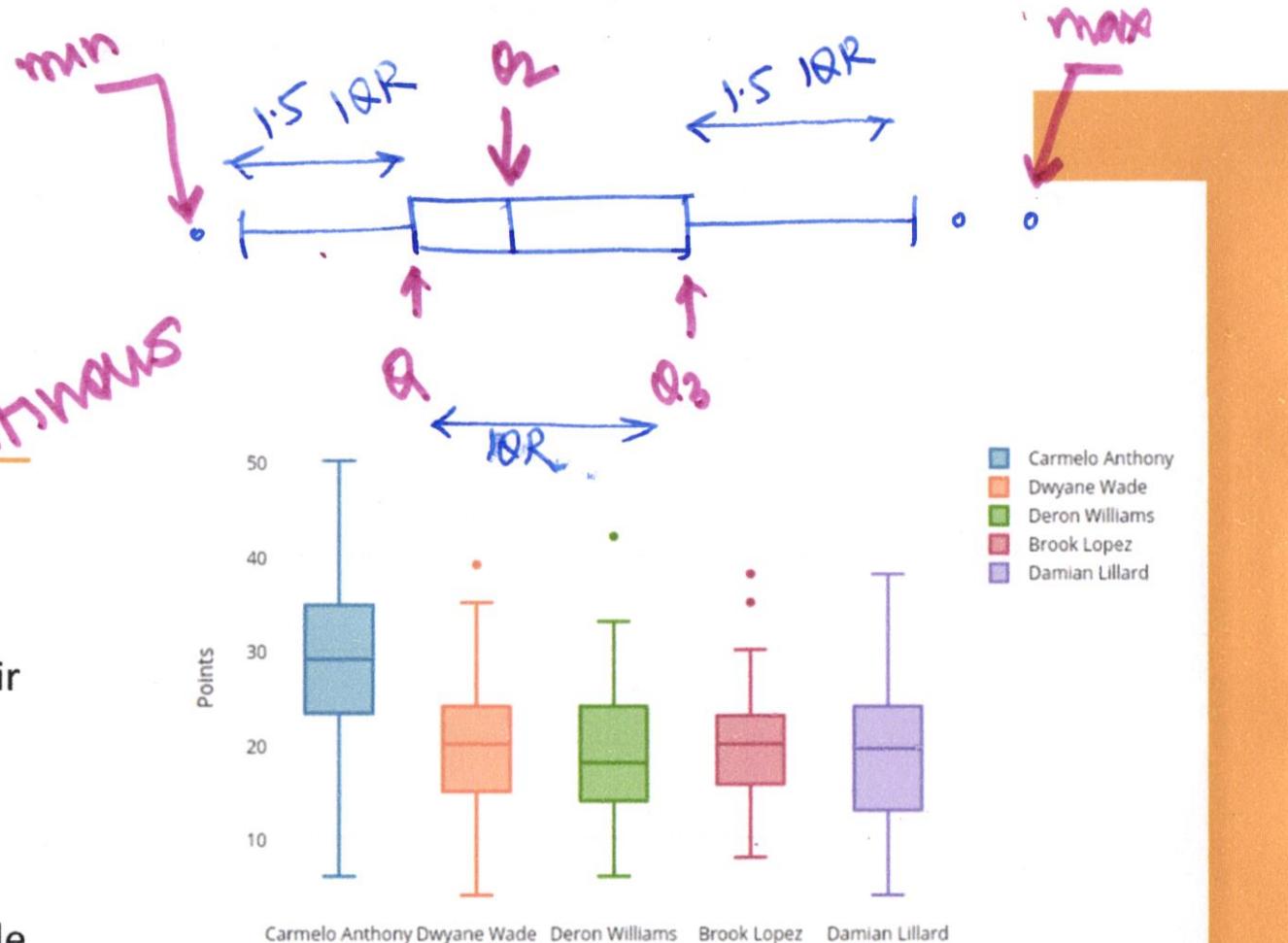


# Box Plot

and whiskers

continues

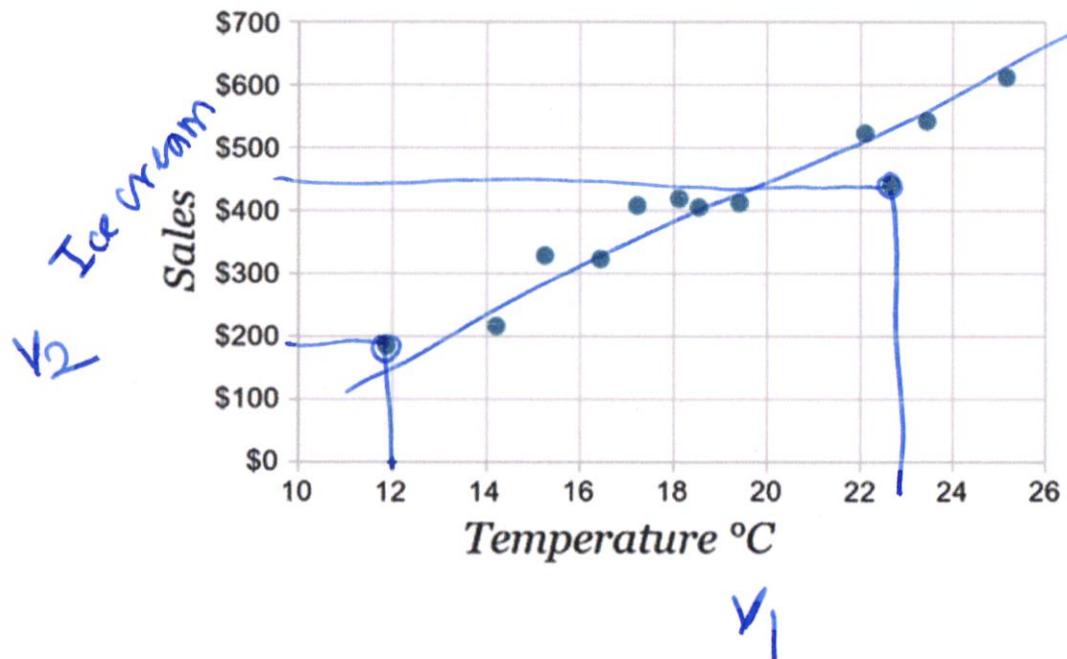
- A method for graphically depicting groups of numerical data through their quartiles.
- Gives five point summary, Minimum, Lower quartile, Median, Upper quartile and Maximum.
- Example: Marketing Budget

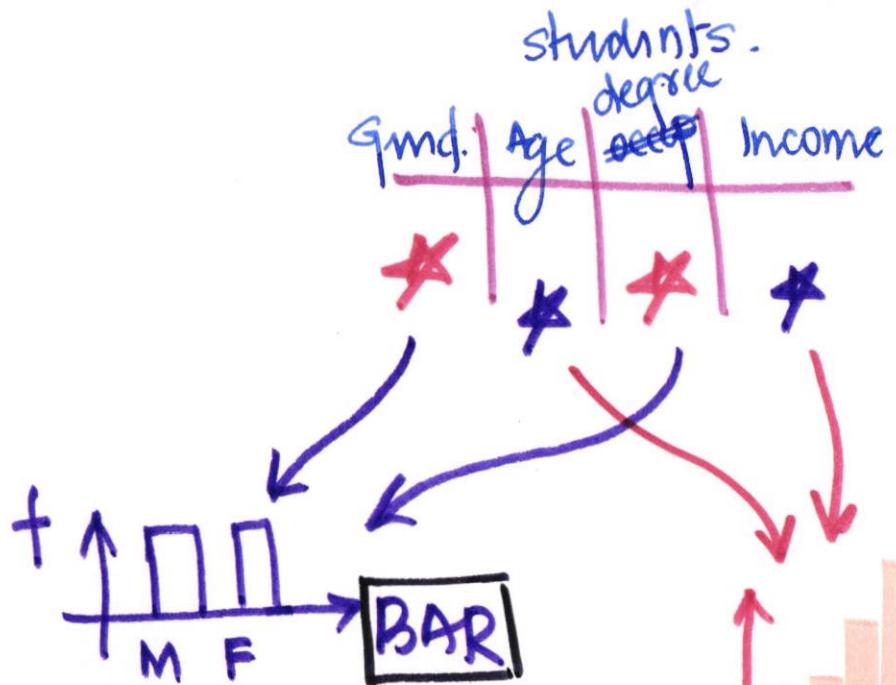


# Scatter Plot

Bivariate plot

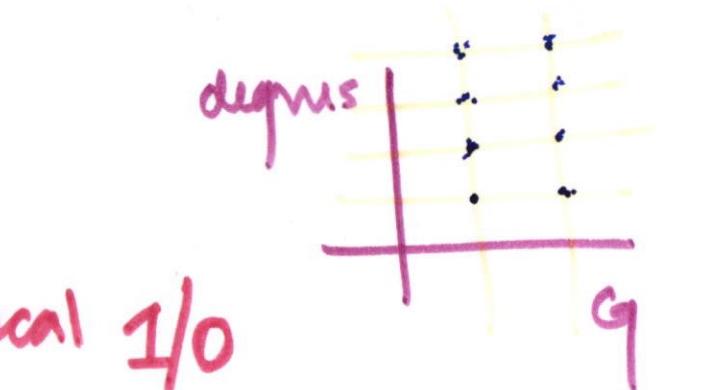
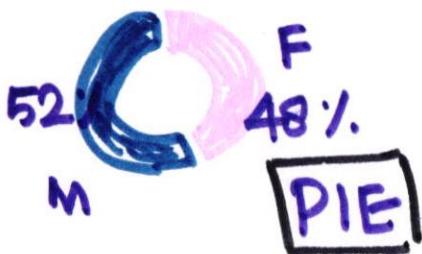
- This graph uses dots to represent values for two different numeric variables.
- Generally, they are used to observe relationships between variables.
- Example: Ice-cream sales





\* - categorical I/O

\* - continuous



**Scatter**  
Bi Variate

