

# Statistics

## Definition -

A branch of applied maths that involves the collection, description, presentation, analysis and interpretation of numerical data. is called statistics.

## Type of Data -

- ① structured data
- ② unstructured data

## Stages of stats -

- ① collection of data

- ② Organizing of data
- ③ Presentation of data
- ④ Analysis of data
- ⑤ Interpretation of data

Type of stats -

### ① Descriptive stats

Available data sample or population on it we perform action like analyzing, describe, summarize. it called descriptive stats.

## ② Inferential stats -

on describe data we perform interpretation like hypothesis testing on the data for example z-test, t-test, F-test, chi-square test. is called inferential statistics.

### ① Descriptive stats.

① univariate Des. stats.

② Bivariate Des. stats.

③ multi variate Des. stats.

## \* Population & Sample

Entire data available is call population.

$$\text{Population} = N$$



$$\text{Sample} = n$$

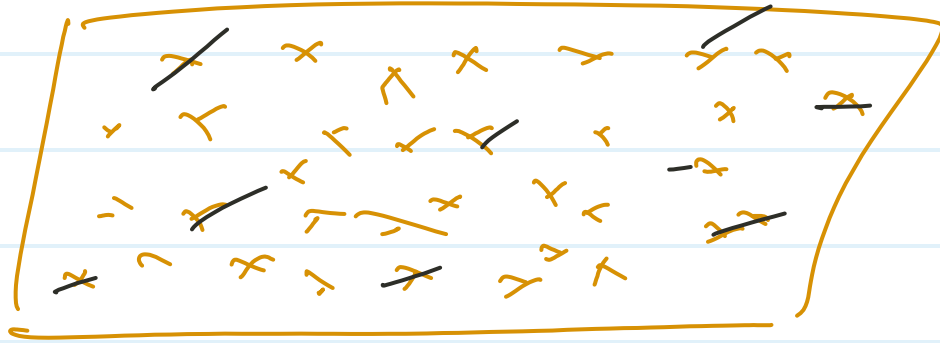
## \* Sampling methods -

Type of sampling method

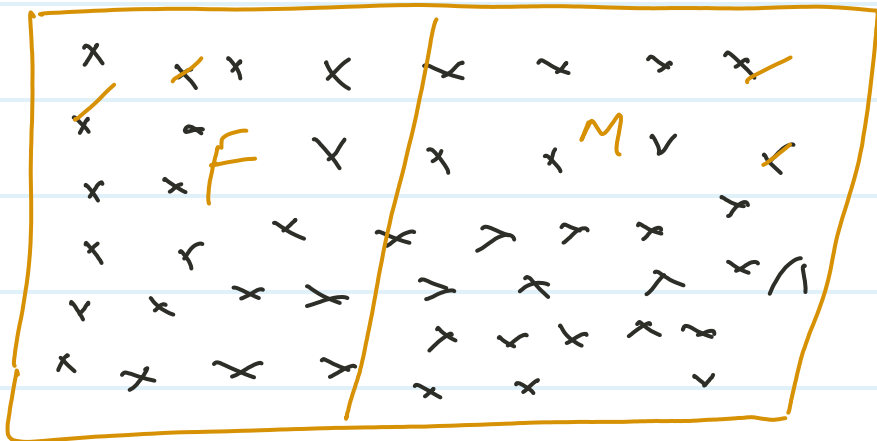
- ① Probability sampling
- ② non-probability sampling.

# Probability sampling

## ① Simple random sampling.-



## ② stratified sampling



> 18

Gender.

### ③ Systematic Sampling.-

$n^{\text{th}}$  selected

population - 10000  
sample - 50

$$= \frac{10000}{50}$$

$$= 200.$$

1<sup>st</sup>, 201<sup>st</sup>, 301<sup>st</sup>,

[✓, |, |, |, |, |, |, |, |, |, |, |]

### ② Non-probability Sampling -

Convenience Sampling -

Research -

# How to calculate sample size?

7

## \* Cochran formula

$$n_0 = \frac{Z^2 p q}{e^2} \quad \text{— infinite population}$$

$$n = \frac{n_0}{1 + \frac{(n_0 - 1)}{N}} \quad \text{— Required population.}$$

$e$  = margin of error  $\leq 1$

$P$  = Population proportion 50%

$$q = 1 - P$$

$Z$  = value from Z-table

C.I.

 $\alpha$ 

90 %	-	0.1
<del>95 %</del>	-	0.05
99 %	-	0.01

z-table

$$\alpha = 0.05$$

$$1 - 0.05 = 0.95$$

$$Z = 1.6 + 0.05$$

$$Z = 1.65$$

$$= \frac{(1.65)^2 \times 0.5 \times (1 - 0.5)}{(0.05)^2}$$

$$= \frac{2.72 \times 0.5 \times 0.5}{0.025}$$

$$n_0 = \frac{0.68}{0.025} = 27.2$$



Required population  $N = 10000$

$$n = \frac{n_0}{1 + \frac{(n_0 - 1)}{N}}$$

$$= \frac{27.2}{1 + \frac{(27.2 - 1)}{10000}}$$

$$\Rightarrow 27.12$$

$$28 / 27$$

if population = 100000

Sample =

# Descriptive stats -

- ① measure of Center tendency
- ② measure of Dispersion or variation
- ③ measure of position
- ④ measure of shape

## \* Measure of Center Tendency

① mean -

[2, 3, 4, 5, 6]

$$\text{mean} = \frac{2+3+4+5+6}{5} = \frac{20}{5}$$

$$= 4$$

population mean =  $\mu$

sample mean =  $\bar{x}$

② median -

I-case  $[6, 2, 4, 5, 3]$

data sort

$[2, 3, 4, 5, 6]$   
 $\uparrow$

median = 4

II-Case

$[2, 3, 5, \underline{7, 8}, 9, 11, 13]$

$$\text{median} = \frac{7+8}{2} \Rightarrow 7.5$$

③ mode -

- (i) uni-modal
- (ii) Bi-modal
- (iii) multi-modal

$$① \quad [2, 3, 4, 5, 5, 6]$$

$$\text{mode} = 5$$

$$② \quad [2, 2, 4, 5, 6, 6]$$

$$\text{mode} = 2, 6$$

$$③ \quad [2, 2, 4, 4, 5, 5, 6, 7, 8]$$

$$\text{mode} = 2, 4, 5$$

→ main application

→ outlier handling

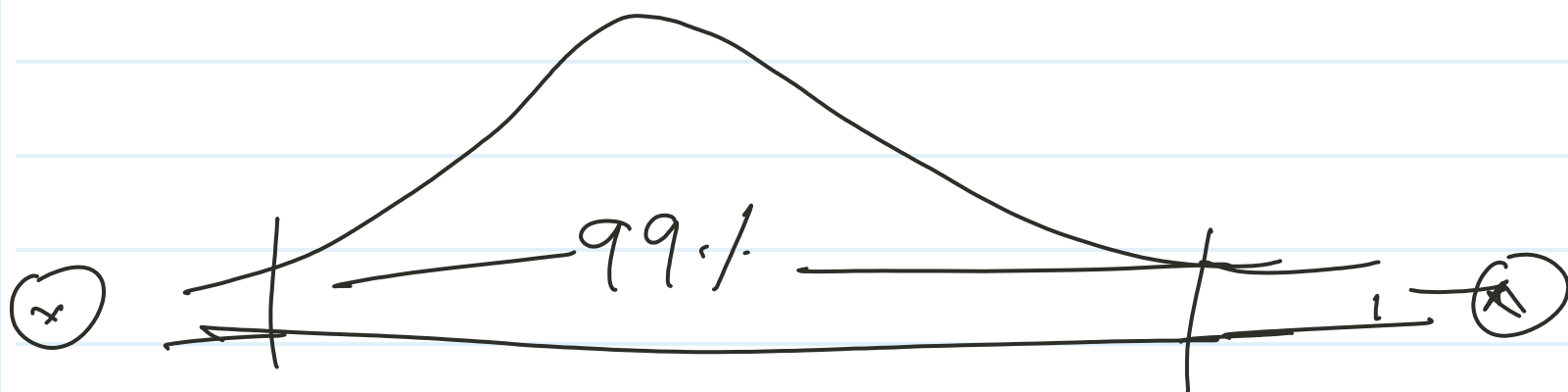
→ missing value handel

X	Y
2	2
3	3
✓	4
7	45 ✓
✓	7
11	5
13	

$$2, 3, 4, 5, 7, 45$$

$$\quad \quad \quad \underline{4.5}$$

$$\underline{2, 3, 4, 5, 7, 4.5}$$



Result

P

F

P

✓

F

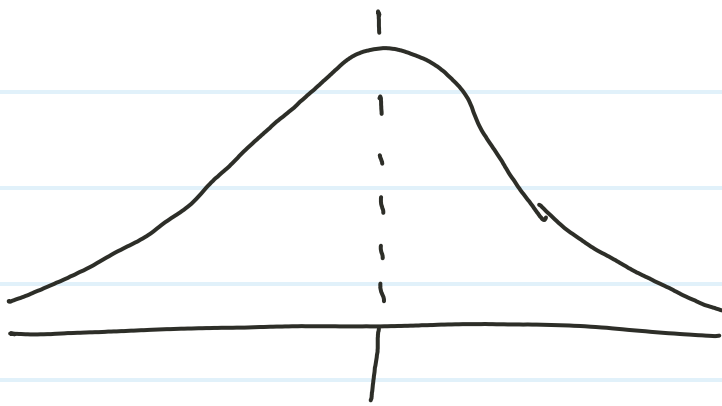
P

✓

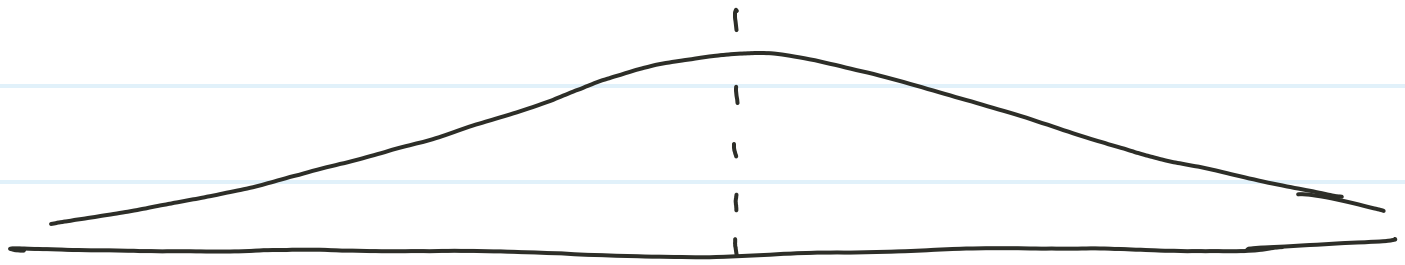
F

P

\* measure of Dispersion or variance



— low variability



High variability

① mean absolute deviation -

The mean absolute devi. of a dataset is the avg. distance b/w each data point and the mean.

$$\Rightarrow \frac{1}{N} \sum_{i=1}^n |x_i - \bar{x}|$$

Data -  $[10, 15, 15, 17, 18, 21]$

$$\text{mean} = \frac{96}{6} = 16$$

$$|10 - 16| = |-6| = 6$$

$$|15 - 16| = |-1| = 1$$

$$|15 - 16| = 1$$

$$|17 - 16| = 1$$

$$|18 - 16| = 2$$

$$|21 - 16| = \frac{5}{16}$$

$$\Rightarrow \frac{16}{6} \Rightarrow 2.67$$

② Variance -

It tells the degree of spread in dataset.  
High variability - Datapoint spread widely

low variability - data point close to mean.

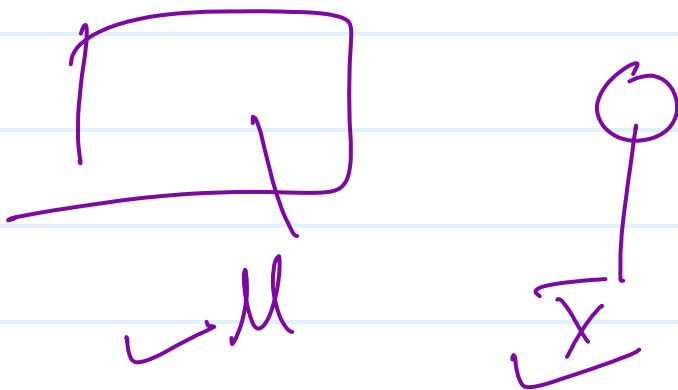
population  $\sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2$

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

Note -  $n-1$  is a degree of freedom

OR

Bessel's correction  
For keep away result from biased.





### ③ standard deviation -

The square root of variance is called st. dev.

The farther the data points from the higher the deviation

$$\text{Pop } \sigma = \sqrt{\frac{1}{N} \sum_{i=1} (x_i - \mu)^2}$$

$$\text{Sample } S = \sqrt{\frac{1}{n-1} \sum_{i=1} (x_i - \bar{x})^2}$$

### ③ Range -

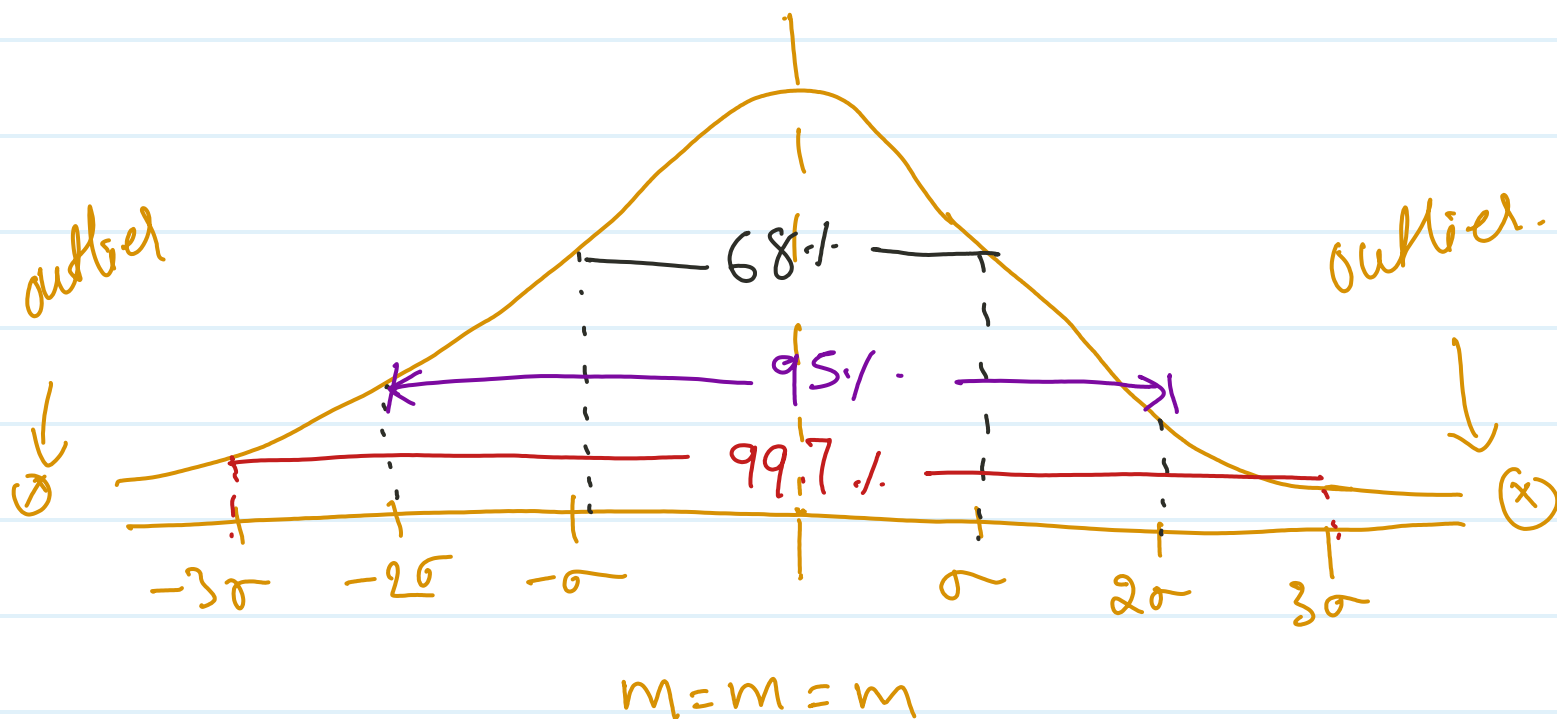
$$[1, 2, 5, 6, 11, 15, 19, 25, 30]$$

$$\text{max} - 30$$

$$\text{min} - 1$$

$$\begin{aligned} \text{Range} &= \text{max} - \text{min} \\ &= 30 - 1 = 29 \end{aligned}$$

# Empirical Rule



## Normal Distribution Curve

68 — 95 — 99.7  
 $\sigma$        $2\sigma$        $3\sigma$

