# Statistics

## Definition —

A branch of applied maths that involves the collection description, presentation, analysis and Interepretation 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 Interepretation of data


Type of stats -

① Descriptive stats

Available data sample or population
on it we perform action like
analysing. discribe, summerize. it called
descriptive stats.

② <u>Inferential stats -</u>

On describe data we perform interepretation 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 multivariate Des. stats.

# 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 -

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

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

$$= 4$$

population mean $= \mu$

sample mean $= \overline{x}$

② median —

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

data   sort

[2, 3, 4, 5, 6]
↑

median = 4

II - case

[ 2, 3, 5, 7, 8, 9, 11, 13]

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

③ mode —

    ① uni-modal
    ② Bi-modal
    ③ multi-modal

① $[2,3,4,5,5,6]$

mode = 5

② $[2,2,4,5,6,6]$

mode = 2, 6

③ $[2,2,4,4,5,5,6,7,8]$

mode = 2, 4, 5

→ main application

→ outlier handeling

→ missing value handel

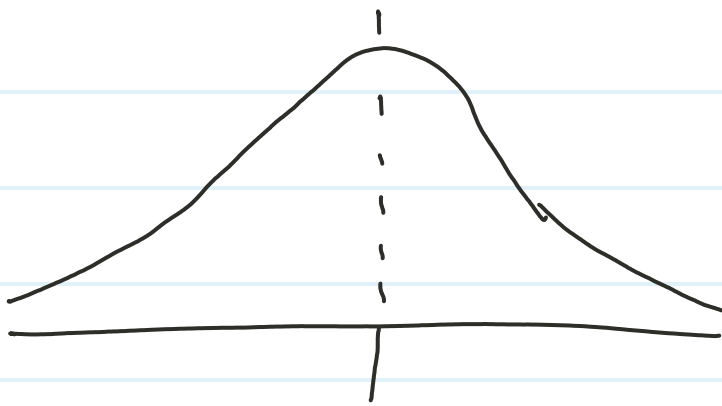| X | Y | |
|---|---|---|
| 2 | 2 | $2,3,4,5,7,45$ |
| 3 | 3 | $\overline{9.5}$ |
| ✓ | 4 | |
| 7 | 45✓ | $2,3,4,5,7,45$ |
| ✓ | 7 | |
| 11 | 5 | |
| 13 | | |

$x$ ⟵ 99% ⟶ $\bar{x}$
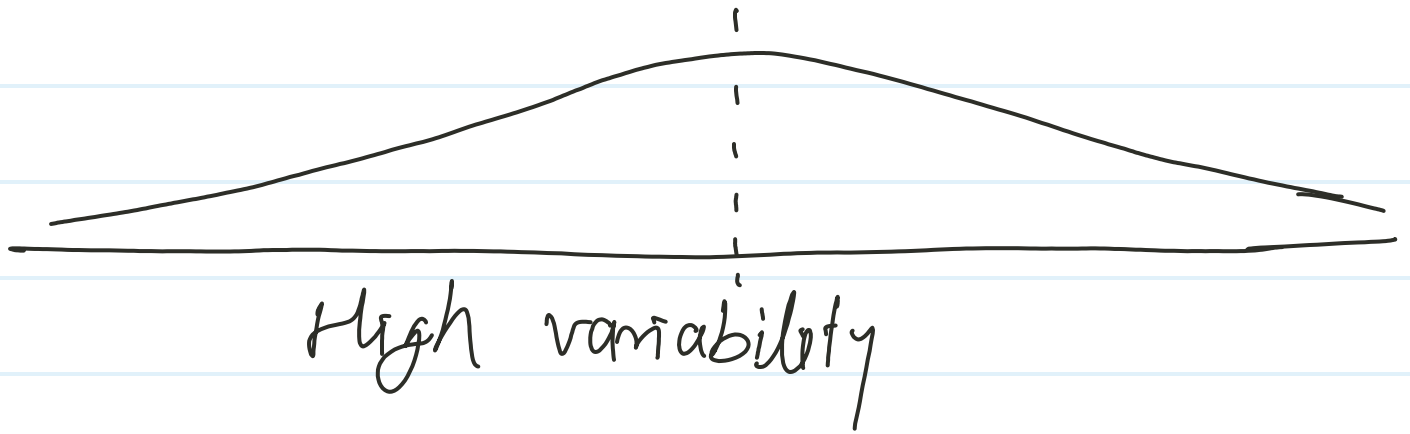
$1$

Result

P

F

P

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.

$$\Rightarrow \quad \frac{1}{N} \sum_{i=1}^{n} |X_i - \overline{X}|$$

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

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

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

$|15-16| \qquad = 1$

$|15-16| \qquad = 1$

$|17-16| \qquad = 1$

$|18-16| \qquad = 2$

$|21-16| \qquad = \dfrac{5}{16}$

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

② Variance -

It tells the degree of spread in dataset. High variability — datapoint spread widely
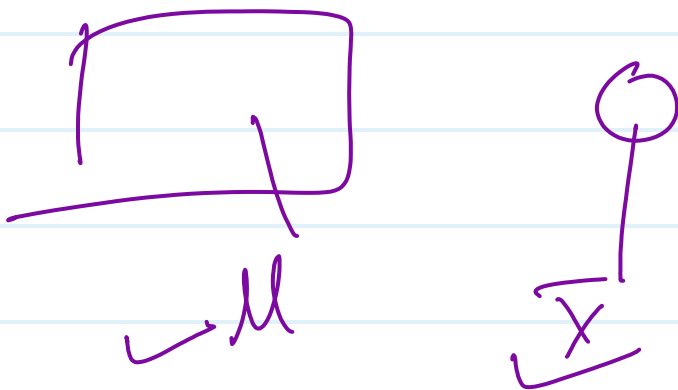
low variability — data point close to mean.

population $\sigma^2 = \dfrac{1}{N} \sum\limits_{i=1}^{n} (X_i - \mu)^2$

Sample $S^2 = \dfrac{1}{n-1} \sum\limits_{i=1}^{n} (X_i - \bar{x})^2$

Note — $n-1$ is a degree of freedom
OR
Besils 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

Pop $\qquad \sigma = \sqrt{\dfrac{1}{N} \sum\limits_{i=1} (X_i - \mu)^2}$

Sample $\qquad S = \sqrt{\dfrac{1}{n-1} \sum\limits_{i=1} (X_i - \bar{X})^2}$

③ Range -

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

max - 30

min - 1

Range = max - min

$\qquad = 30 - 1 \qquad = 29$

# Emperical Rule



outlier

outlier

68%

95%

99.7%

$-3\sigma$    $-2\sigma$    $-\sigma$      $\sigma$    $2\sigma$    $3\sigma$

$m = m = m$