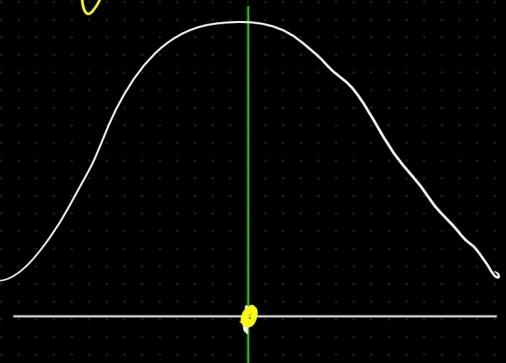


Measure of Central Tendency

① Mean (Average)

② Median

③ Mode



Population
↓ (N)

Sample
(n)

Mean: By adding up all the components
and dividing the number of
components

$$= \sum_{i=1}^N \frac{x_i}{N} \rightarrow \begin{matrix} \text{population} \\ \text{Mean} \end{matrix}$$

$$= \frac{x_1 + x_2 + x_3 + x_4 + \dots + x_n}{N}$$

Population Mean (μ)

$$\mu = \sum_{i=1}^N \frac{x_i}{N}$$

Sample Mean (\bar{x})

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

$$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$

$$= \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= \frac{32}{10} = 3.2 = \text{Mean}$$

(2) Median : The median is the middle value of a dataset.

$$X = \{4, 6, 5, 9, 10, 2, 13\}$$

Steps : (1) Sort the Data

$$\{1, 2, 4, 5, 6, 9, 10\}$$

(2) No of elements \rightarrow (7)

$$\text{Count} = 7$$

3) If count is odd

$$M = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ term}$$

$$M = \left(\frac{7+1}{2} \right)^{\text{th}} \Rightarrow 4^{\text{th}}$$

$$\boxed{\text{Median} = 5}$$

$$Y = \{1, 2, 4, 6, 10, 12\}$$

$$\text{Count} = 6$$

$$\text{Count} = \text{even}$$

$$M = \frac{\left(\frac{N}{2} \right)^{\text{th}} + \left(\frac{N}{2} + 1 \right)^{\text{th}}}{2}$$

$$= \frac{3^{\text{th}} + 4^{\text{th}}}{2} = \frac{4+6}{2} = \frac{10}{2}$$

$$\boxed{\text{Median} = 5}$$

↓ Outlier

$$X = \{1, 2, 3, 4, 5, 100\}$$

$$X = \{1, 2, 3, 4, 5\}$$

$$\bar{x} = \frac{1 + 2 + 3 + 4 + 5}{5}$$

$$\bar{x} = \frac{15}{5} = 3$$

$$X = \{1, 2, 3, 4, 5, 100\}$$

$$\bar{x} = \frac{1 + 2 + 3 + 4 + 5 + 100}{6}$$

$$= \frac{115}{6} \approx 19.$$

(16)

3 → 19
↓
outliers

$$\text{Median} = \frac{3 + 4}{2} = \frac{7}{2} = 3.5$$

Note: Median is used to find the central tendency when outliers are present.

③ Mode: The mode is the value that occurs most often.
[Frequency Maximum]

$\{10, 20, 30, 10, 10, 40\}$

[Mode = 10]

NAN → Missing Values → handle
↓

$E^D A$
=

Age

20

10

30

40



—

50

—



Mean / Median
=

Class

A

B

C

—

—

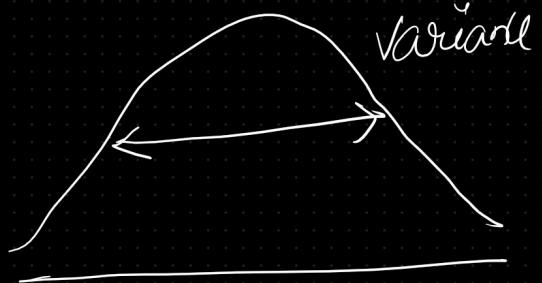


Mode
=

Measure of Dispersion [spread of the data]

1) Variance

2) Standard deviation



Population Variance

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

x_i → data points

μ → population mean

N → population size

Sample Variance

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

x_i → data points

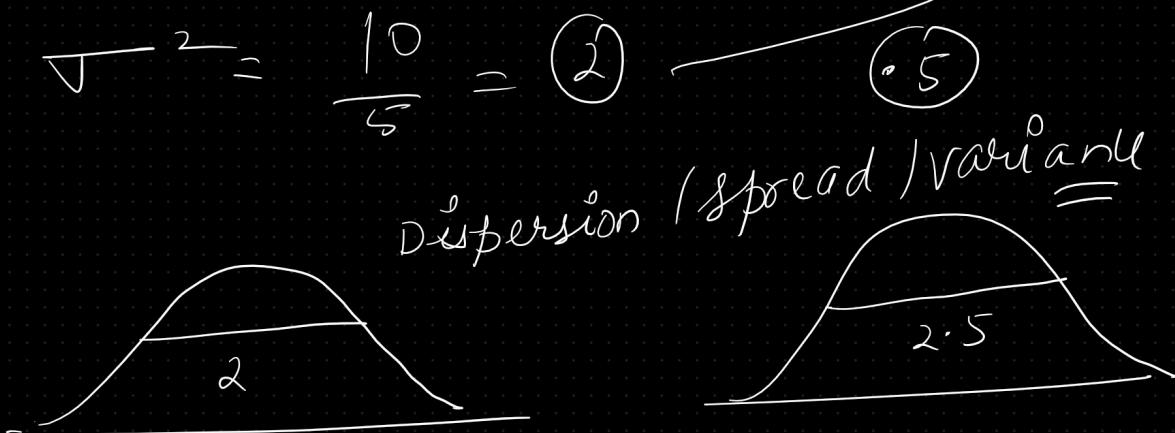
\bar{x} → sample mean

$n \rightarrow$ Sample size

$$X = \{1, 2, 3, 4, 5\}$$

$$\begin{array}{ccc}
 x_i & \bar{x} & (x_i - \bar{x})^2 \\
 1 & 3 & 4 \\
 2 & 3 & 1 \\
 3 & 3 & 0 \\
 4 & 3 & 1 \\
 5 & 3 & 4 \\
 \hline
 \bar{x} = 3 & \overline{(10)} & = \frac{10}{5-1} = \frac{10}{4} = 2.5
 \end{array}$$

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



Standard deviation

population std

$$\sigma = \sqrt{\text{Variance}}$$

Sample std

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

↓

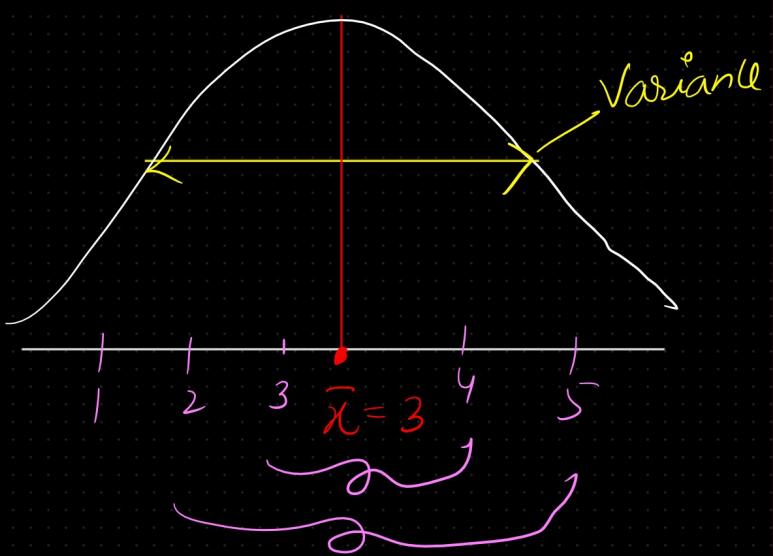
Sample Variance

$$X = \{1, 2, 3, 4, 5\}$$

$$\bar{x} = 3$$

→ How much individual data points differ from the mean

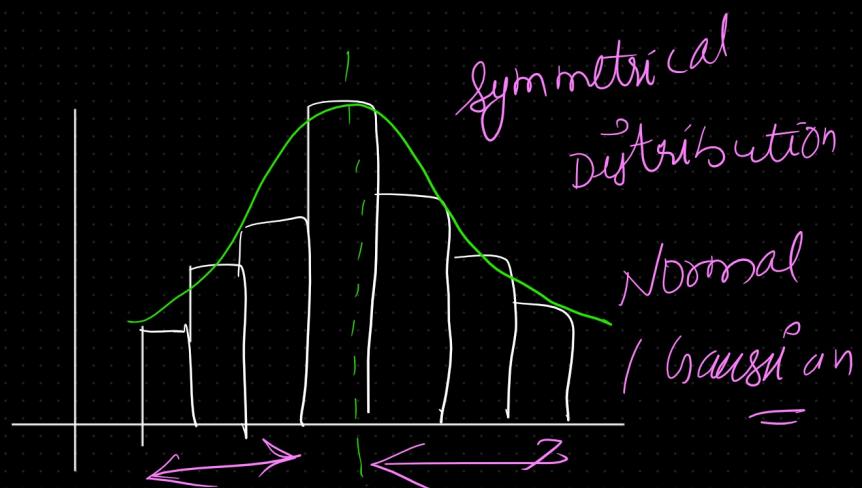
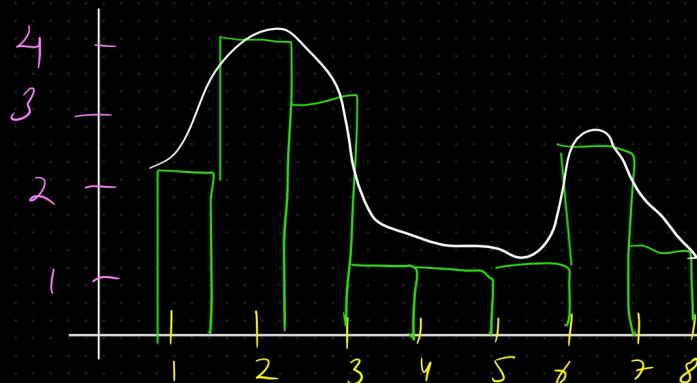
→ It provides a way to express how spread out of the values in a dataset.



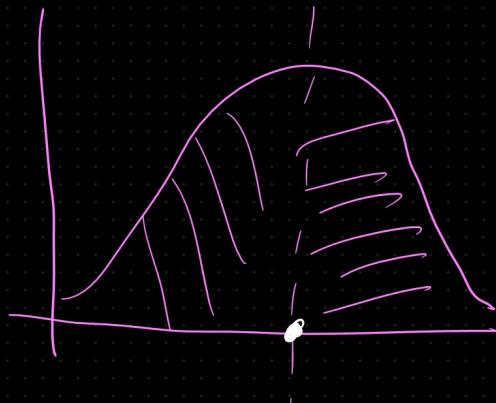
Histogram And Skewness

[Frequency]

data = {1, 1, 2, 3, 3, 3, 2, 2, 2, 4, 5, 6, 8, 16, 6, 8}



→ NO Skewness



Mean, mode,
median



$$Q_2 - Q_1 \approx Q_3 - Q_2$$

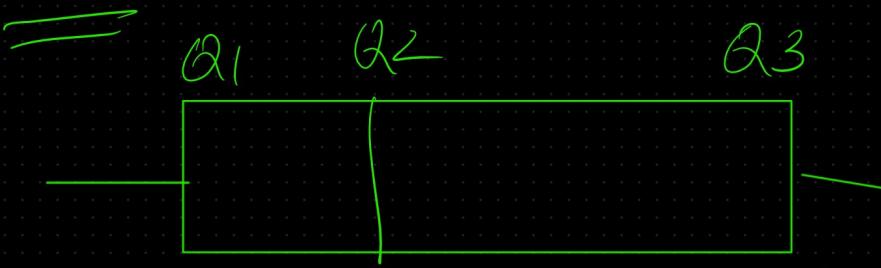
mean = mode = median

② Right Skewed



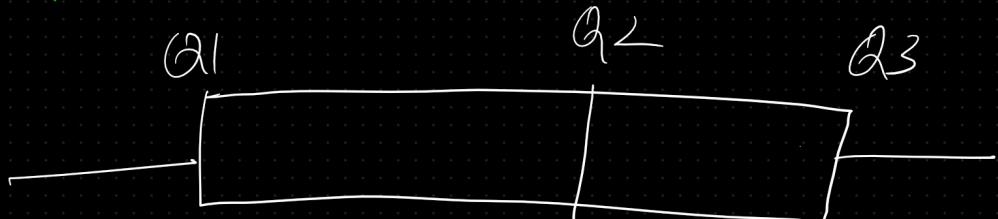
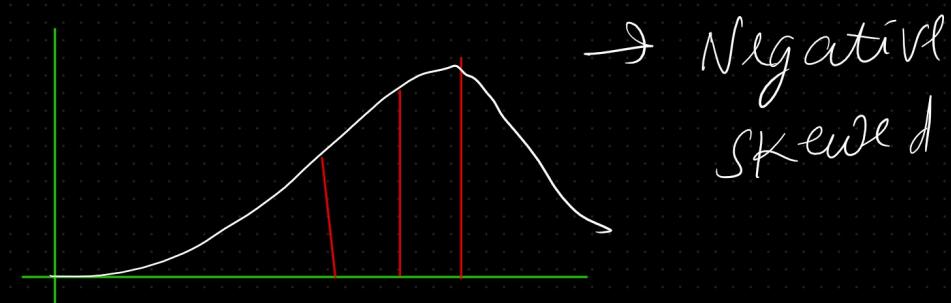
mean \geq median \geq mode

Box Plot



③

Left Skewed Distribution



$$Q_2 - Q_1 \geq Q_3 - Q_2$$

mean \leq median \leq mode

Covariance and Correlation

X	Y	[Relationship b/w X and Y]
2	4	X↑ Y↑
5	6	X↓ Y↓
7	8	X↓ Y↑
10	12	X↓ Y↓

Covariance

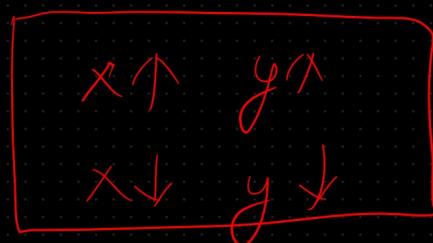
$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

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

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

= $\text{Cov}(x, x) \Rightarrow \text{spread}$

$\text{Cov}(x, y)$



+ve Cov

$\text{Cov}(x, y)$



-ve Cov

Coefficient

[-1 to 1]

(i) Pearson Correlation Coeff-

$$\rho_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y}$$

a)

$$+1 \swarrow \text{Cov}(x, y)$$

b)

$$-1 \curvearrowleft \text{Cov}(x, y)$$

The more the value towards +1

the more +ve correlated is (x, y)

The more the values towards -)
the more -ve correlated (x, y)

(2) Spearman Rank Correlation
[-1 to 1]

$$\rho_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$$