

general
Things

Algebraic Variable = $x + 5 = 10$

$$x = 5$$

Random Variable = x

Throwing a dice



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

List Variable = $x = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$

List

Statistical Measure's

Measure of Central Tendency

— Measure the Central Value of Dataset.

Mean

Numerical

The mean is sum of all observation divided by total no. of observation.

$$\text{data} = [1, 2, 3, 4, 5]$$

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5$

$$\text{Mean} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

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

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

Population Mean	Sample Mean
$\mu = \frac{\sum_{i=1}^N x_i}{N}$ <p>N = number of items in the population</p>	$\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$ <p>n = number of items in the sample</p>

Median

Numerical

The median is the middle value in the dataset, when the data is arranged in order.

$$\text{odd} = \left\lceil \frac{n+1}{2} \right\rceil$$

$$\text{even} = \frac{\left\lceil \frac{n}{2} \right\rceil + \left\lceil \frac{n}{2} + 1 \right\rceil}{2}$$

$$\text{data} = [5, 3, 4, 1, 2]$$

$$\text{data} = [1, 2, 3, 4, 5]$$

x_1, x_2, x_3, x_4, x_5

arrange the data

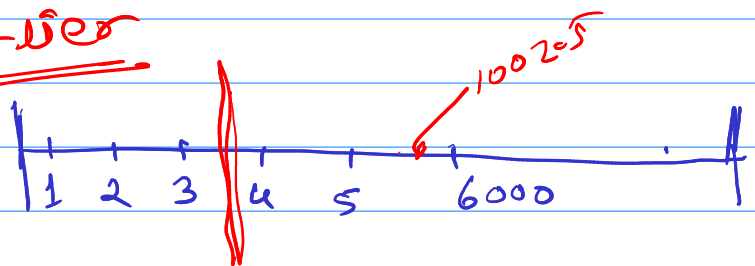
Why median?

Median is Robust to outliers

$$\text{ex} = \text{data} = [1, 2, 3, 4, 5, 6000]$$

$$\bar{\mu} = \frac{1+2+3+4+5+6000}{6}$$

outlier



$$= \frac{6015}{6}$$

$$\bar{\mu} = 1002.5$$

median = [1, 2, 3, 4, 5, 6] — even

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

median = 3.5

— / —

Age = [20, 21, 22, 23, 21, 25, 21, 26, 20, 21, ...]
(60 Student)

20, 20
first

21, 25, 25, 20
Last

Mode = The most frequent value in Dataset.
Categorical variable.

Gender =

F
M
F
F
M

F = 4 = There are two mode
m = 4

mode = M / F

Null

Blank

F
F
M
F

F = 6
m = 4

mode = F

Weighted Mean

- The sum of all the product with its weight, divided by sum of all weight.

AI = Algorithm	name	Accuracy	output → House price prediction
✓ LR		0.2	10L
✓ RF		0.3	15L
✓ Xgboost		0.1	8L
		weights	

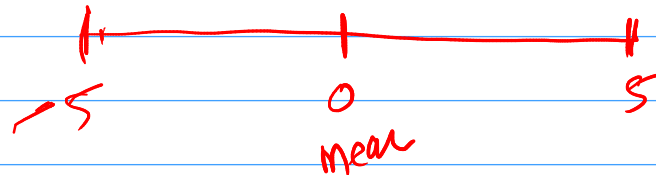
$$= \frac{10L \times 0.2 + 15L \times 0.3 + 8L \times 0.1}{0.2 + 0.3 + 0.1}$$

= weighted mean

output	expenses	height
	-	Bus
	-	Car
	-	Bike
	-	Train

Measure of Dispersion \rightarrow Spread. पसरना

A measure of dispersion is a statistical measure that describes the spread or variability of a dataset. It provides information about how the data is distributed around the central tendency (mean, median or mode) of the dataset.



Range

Range: The range is the difference between the maximum and minimum values in the dataset. It is a simple measure of dispersion that is easy to calculate but can be affected by outliers.

Variance

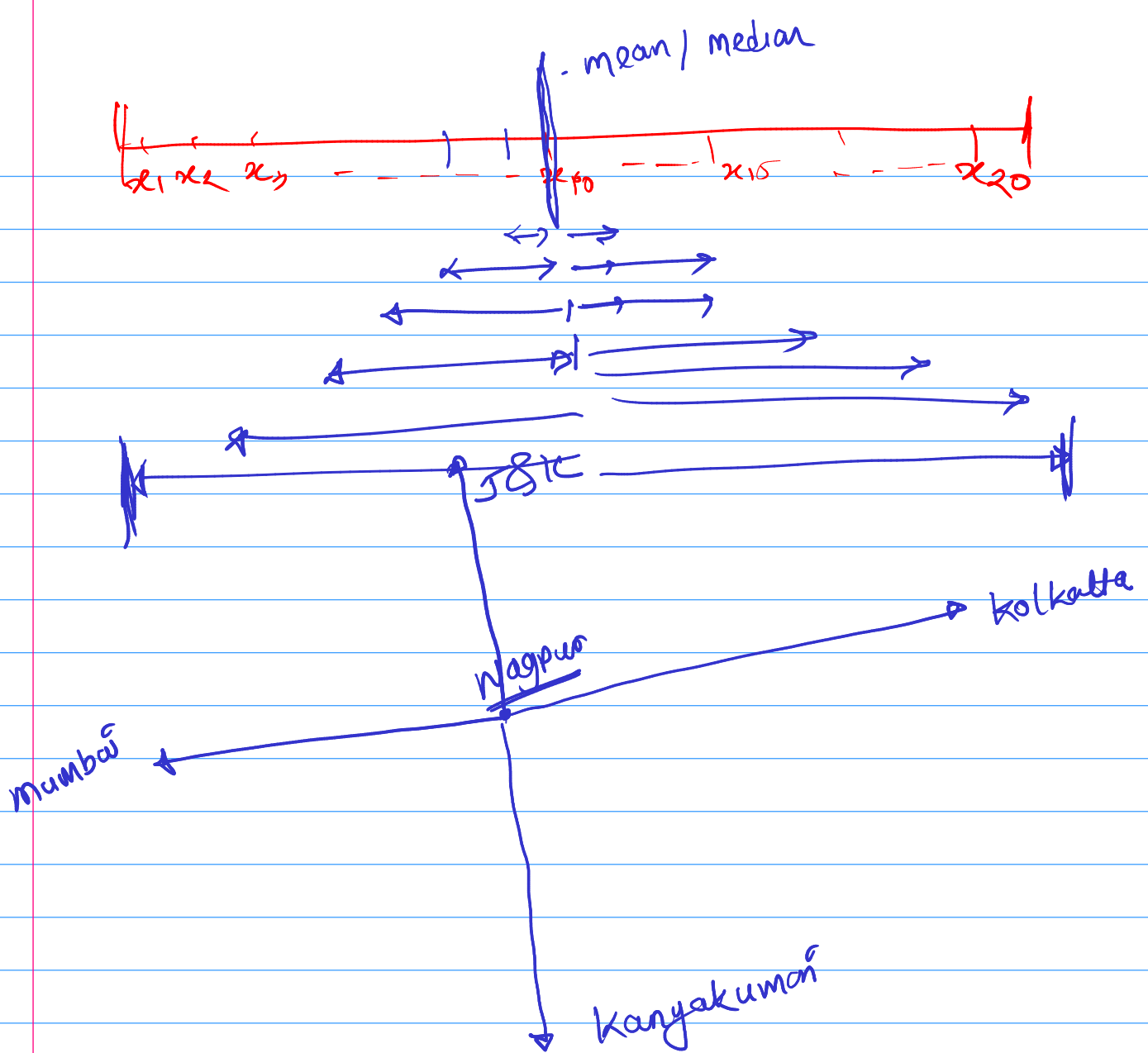
Variance: The variance is the average of the squared differences between each data point and the mean. It measures the average distance of each data point from the mean and is useful in comparing the dispersion of datasets with different means.

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

why?

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

\downarrow sigma.



km \rightarrow km²

lakh \rightarrow lakh²

sq. feet \rightarrow sq. feet²

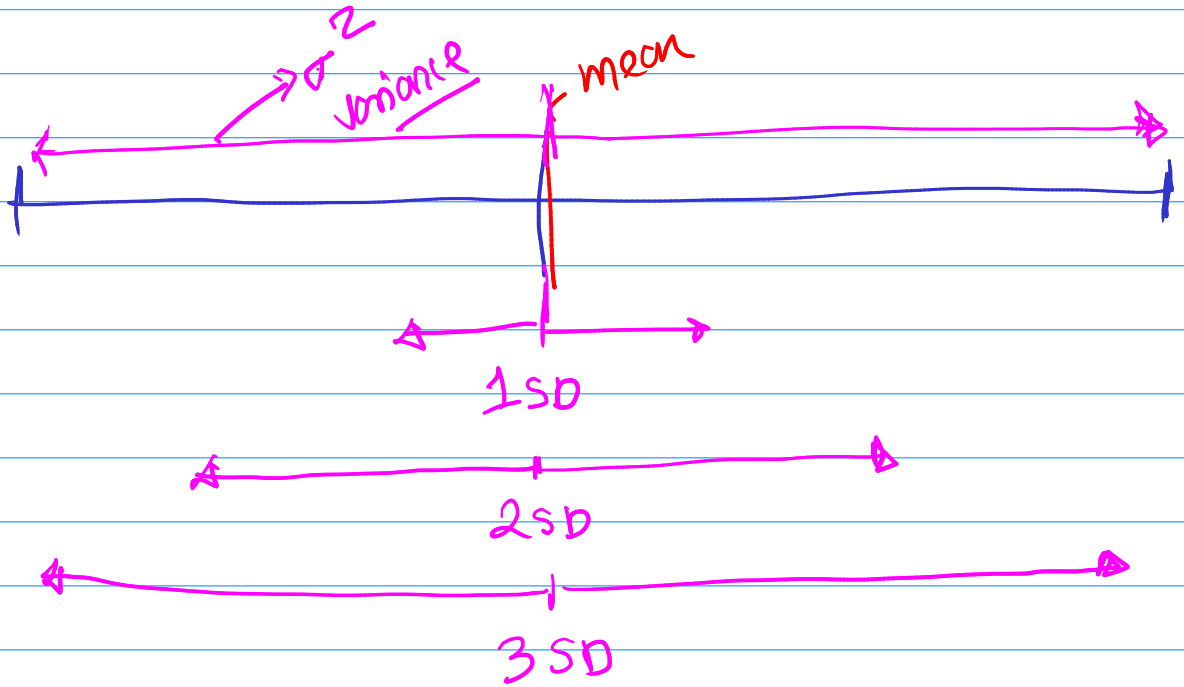
3. Standard deviation

- Standard Deviation:** The standard deviation is the square root of the variance. It is a widely used measure of dispersion that is useful in describing the shape of a distribution.

(SI unit is same)

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

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



why $(n-1)$ in sample variance



$\sigma^2 \approx s^2$ when we divide by $(n-1)$