

# Statistics

⇒ Statistics is a field that deals with collection, organization, analysis, interpretation and presentation of data



Decision making

## Statistics

Descriptive

(It consists of organizing and summarizing of data)

1. Measure of Central Tendency

↳ Mean, Median, Mode

2. Measure of Dispersion

↳ Variance, Standard Deviation

Inferential

Collect Data

↓ using some exp.  
{ \* Z-Test, t-Test }

Conclusion or Inferences

↓ Other data

Population data.

Population =  $N$

Sample =  $n$ .

## Measure of Central Tendency

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

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

## Measure of Dispersion

$$\text{Age 1} = \{2, 2, 4, 4\}$$

$$\text{Age 2} = \{1, 1, 5, 5\}$$

$$\mu = 3$$

$$\mu = 3. \quad \text{Spread is more?}$$

### 1. Variance (Population Data)

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

$x_i$	$\mu$	$(\mu - x_i)^2$
2	3	1
2	3	1
4	3	1
4	3	1
		<hr/>
		4

$$N = 4$$

$$\sigma^2 = 1$$

### (Sample Data)

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

Overcome underestimating of True Population Variance

$x_i$	$\mu$	$(x_i - \mu)^2$
1	3	4
1	3	4
5	3	4
5	3	4
		<hr/>
		16

$$N = 4$$

$$\sigma^2 = 4$$



## Bessel's Correction

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

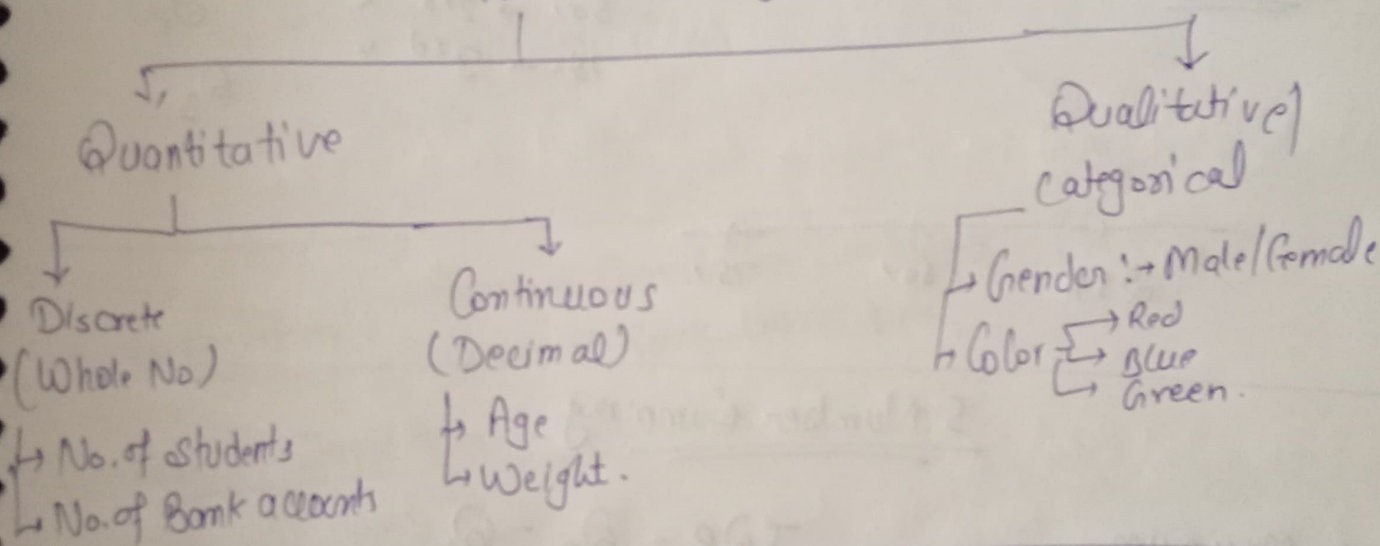
if divided by  $n$  then  
 $\bar{x} < \mu$   
 $s^2 < \sigma^2$  underestimate the true population variance

Standard Deviation  
→ ~~Variance~~ denotes the distance of a data point from Mean.

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

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

Variable (entity which can have any value)



## Random Variable

$X \rightarrow$  function  $\rightarrow$  value  
 $\Downarrow$   
Process / Experiments ( $y = 5x + 2$ )

$X = \left\{ \begin{matrix} 0 & H \\ 1 & T \end{matrix} \right\}$  Tossing of Coin

# Histograms

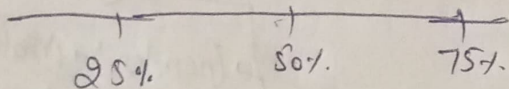
⇒ Percentiles : → is a value below which certain percentage of observation lies.

$$\text{Percentile of } x = \frac{\text{no. of values below } x}{n} \times 100$$

$$\Rightarrow \text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

⇒ Quartiles

25% → 1<sup>st</sup> Quartile  $Q_1$   
50% → 2<sup>nd</sup> "  $Q_2$   
75% → 3<sup>rd</sup> "  $Q_3$



## 5 Number Summary

1. Minimum
2. 1<sup>st</sup> Quartile
3. Median
4. 3<sup>rd</sup> Quartile
5. Maximum

$$IQR = Q_3 - Q_1$$

$$\text{Lower fence} = Q_1 - 1.5(IQR)$$

$$\text{Higher fence} = Q_3 + 1.5(IQR)$$



# Covariance & Correlation (C&C)

C&C are two statistical measures used to determine the relationship b/w two variables. Both are used to understand how changes in one variable are associated with changes in another variable.

⇒ Covariance → is a measure how much two random variables change together. If the variable tends to inc. and dec. together, the covariance is +ve. If one tends to inc. when the other dec. the covariance is -ve.

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

$$\text{Cov}(x, x) = \text{Var}(x) = \sigma^2$$

→ Correlation → Pearson Correlation Coefficient  
Spearman Rank Correlation

① Pearson Correlation Coefficient → [-1 to 1]

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

① The more the value towards +1 the more +ve correlated x & y  
" " " " -1 " " -ve " " "

②

Spearman Rank Correlation

$\rho = -1$ ,  $\rho = -1$  to 0,  $\rho = 0$ ,  $0 < \rho < 1$ ,  $\rho = +1$ .

Non-Linear data

# Probability

⇒ Mutual Exclusive Event : → two events can't occur same time.

⇒  $P(H \text{ or } T) = P(H) + P(T)$  (Additive Rule for MEE)

⇒ Non-mutual exclusive event : ↓

\* Taking out card from deck

$$\begin{aligned} P(K \text{ or } Q) &= P(K) + P(Q) - P(K \text{ and } Q) \\ &= \frac{4}{52} + \frac{12}{52} - \frac{1}{52} = \frac{18}{52} = \frac{4}{13} \end{aligned}$$

⇒ Multiplication Rule

⇒ Independent Event : → 2 events don't affect each other.

② •  $P(H \text{ and } T) = P(H) * P(T)$   
 $= \frac{1}{2} * \frac{1}{2} = \frac{1}{4}$

⇒ Dependent Event : ↓

\* Take a King from the deck and then the Queen card from the deck.

$P(K \text{ and } Q) = P(K) * P(Q|K)$  → Conditional Probability