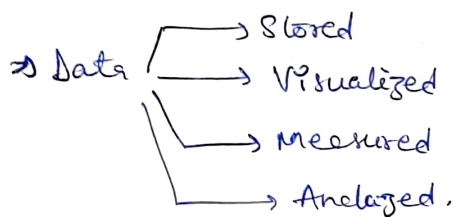


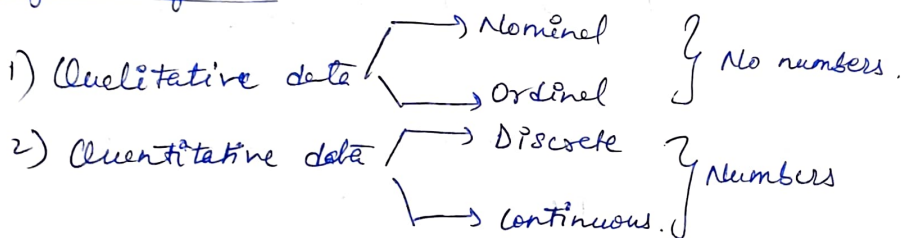
# PART B → STATISTICS and PROBABILITY.

## ① → About data

→ Refers to facts and statistics collected together for reference (or) analysis.



## ② Categories of data



### Nominal

→ No order

Eg: → Male/Female

### Ordinal

→ Basic order

Eg: → Good / Avg / Bad

### Discrete

→ Also called categorical data

Eg: → No: of students in class.

→ Holds only finite number of possible value.

### Continuous

→ Data can hold infinite number of values.

Eg:- Weight of person (50, 50.1, 50.02).

### Note:

Discrete variable → a = Spam / Not spam;

Continuous variable → a = weight;

(Dependent & independent variable also there).

### ③ What is Statistics?

→ Area of applied mathematics concerned with the data collection, analysis, interpretation and presentation.

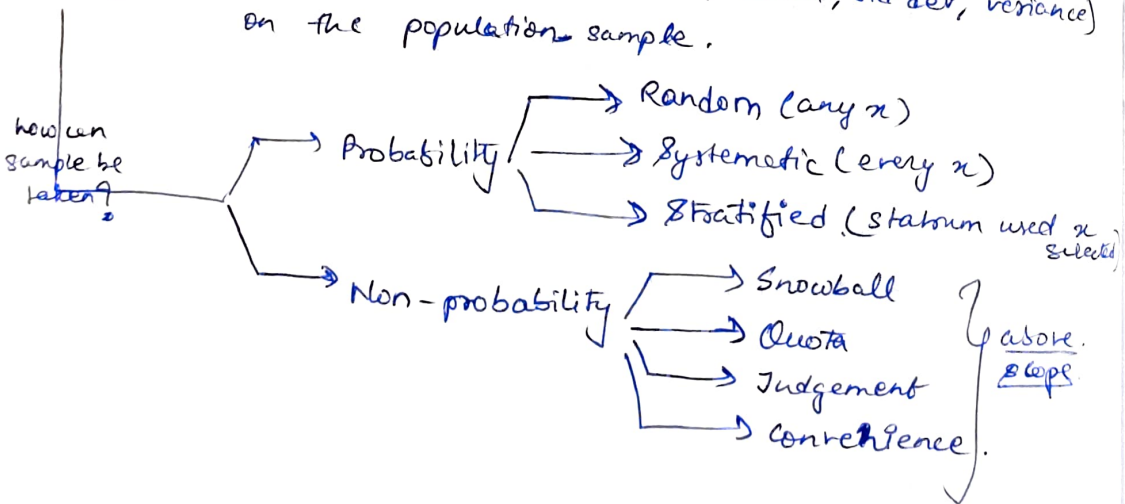
Eg:- i) Data report analysis preparation for a company to identify areas to improve business in.

ii) whether to take a bet or not? If yes, which side should you be on.

### Basic terminologies

→ population, sample, Sampling (taking a sample & inferring).

Sampling → To do infer statistics (mean, median, std dev, variance) on the population sample.



### Types of Statistics

Descriptive → Uses data to provide descriptions of the population, either through numerical calculations (or) graphs (or) tables.  
(Describes data)

Inferential → Inferences and predictions about a population.  
(Infers from data).  
∴ Derives conclusions.

## ④ Descriptive Statistics

→ Describes (or) summarizes data.

→ Des Stats (categories)   
  $\left\{ \begin{array}{l} \rightarrow \text{Measures of Central Tendency.} \\ \rightarrow \text{Measures of Variability/Spread.} \end{array} \right.$

### Measures of Central Tendency

- ↳ Mean (Sum/n)
- ↳ Median (Order & middle value)
- ↳ Mode (most repeated value).

### Measure of Spread

- ↳ Range
- ↳ Inter Quartile Range
- ↳ Variance
- ↳ Std deviation.

→ (i) Range  $\rightarrow \text{Max} - \text{Min}$

→ (ii) IQR  $\rightarrow$  Quartiles tell us about spread of dataset by breaking the dataset into quarters.  
LIKE median breaks it in half.

$$\rightarrow \text{IQR} = Q_3 - Q_1 \quad (\text{Each quarter is } 25\%)$$

→ (iii) Variance  $\rightarrow$  How much a random var differs from it's expected value. Sq of deviation.

(iv) Note: 
$$\rightarrow \frac{\sum (x_i - \bar{x})^2}{n} \quad \left[ \begin{array}{l} \bar{x} \rightarrow \text{mean} \\ x_i \rightarrow \text{element (any)} \\ n \rightarrow \text{total data points} \end{array} \right]$$

Deviation  $\rightarrow$  Difference between each element from it's mean.  
[Not std deviation]

$$\rightarrow \text{ } (x_i - \mu) \text{ (or) } (x_i - \bar{x})$$

$\rightarrow$  Population and Sample Variance.

(n)	(n-1)
$\sigma^2$	$s^2$

(v) Std deviation  $\rightarrow$  Dispersion of a set of data from it's mean.

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

## (VI) Information Gain & Entropy

→ Used in random forest and decision tree, to decide what will be root node.

→ Entropy → Measure of uncertainty present in data.

$$\rightarrow H(S) = - \sum_{i=1}^N p_i \log_2(p_i) \quad \left[ \begin{array}{l} p_i \rightarrow \text{event probability} \\ S \rightarrow \text{set of all instances} \end{array} \right]$$

Information Gain →  $\text{Gain}(A, S) = H(S) - H(A, S)$

$$\left[ \begin{array}{l} H(A, S) \rightarrow \text{entropy of attribute } A \\ H(S) \rightarrow \text{entropy of dataset } S \end{array} \right]$$

→ How much 'Information' a particular feature/variable gives about final outcome.

→ Entropy  $\propto \frac{1}{\text{Significance of variable in decision.}}$

$\therefore I_G \propto$  how much info this var will provide.

## (VII) Confusion Matrix

→ A table that is often used to describe performance of a classification model, on a set of test data.

∴ It represents a tabular representation of Actual vs predicted values.

$$\rightarrow \frac{(\text{True +ve}) + (\text{True -ve})}{(\text{True +ve}) + (\text{True -ve}) + (\text{False +ve}) + (\text{False -ve})}$$

	Yes	No	Actual
Yes	TP	FP	
No	FN	TN	
Can Predicted			