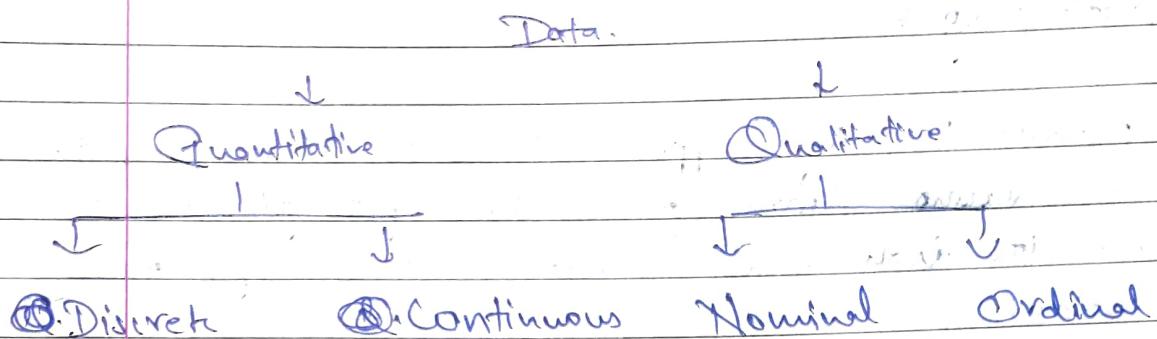


## Statistics.



- Deals with characteristics & descriptors that cannot be easily measured but can be observed effectively.

Nominal → Data with no inherent order  
 Ordinal → Data with an ordered series.

### Quantitative Data

- Deals with nos. of things that could be measured objectively.

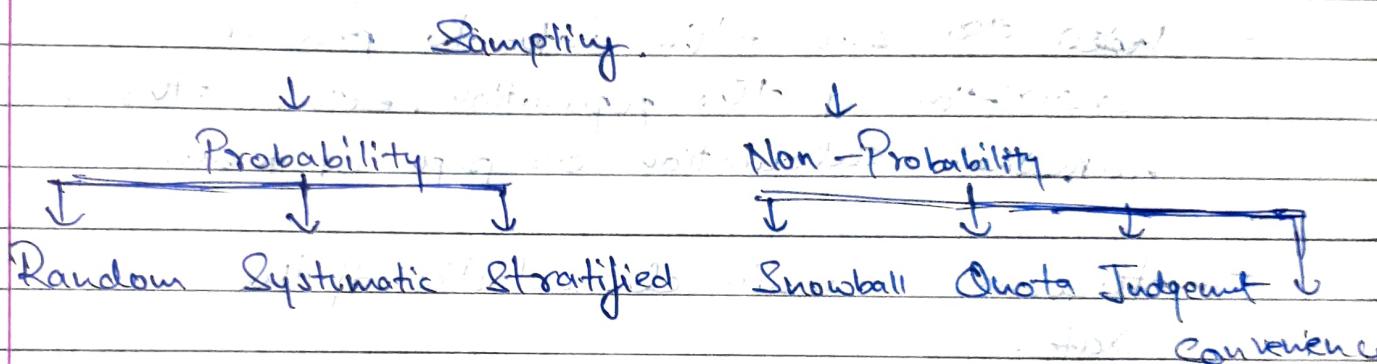
Discrete = AKA Categorical data, holds finite no. of possible values. (No. of students in class)

Continuous = Data can hold infinite no. of possible values.  
 (Weighting scale.)

Population - Collection of or set of individuals or ~~and~~ objects or events whose properties are to be analyzed.

Sample - A subset of population

Sampling - Is a technique that deals with the selection of individual observations within a population. Used to infer statistical knowledge about population.



Random Sampling : A sampling technique in which samples from a large population are chosen using the theory of probability.

Random Sampling : Each member of the population has equal chance to be selected in the sample.

Systematic Sampling : Every nth record is chosen to be a part of the sample.

Stratified Sampling : A stratum is used to form sample. A stratum is a subset of the population that shares at least one common characteristic.

Random Sampling is used to select a sufficient no. of subjects from each stratum.

## Types of Statistics

### ① Descriptive Statistics

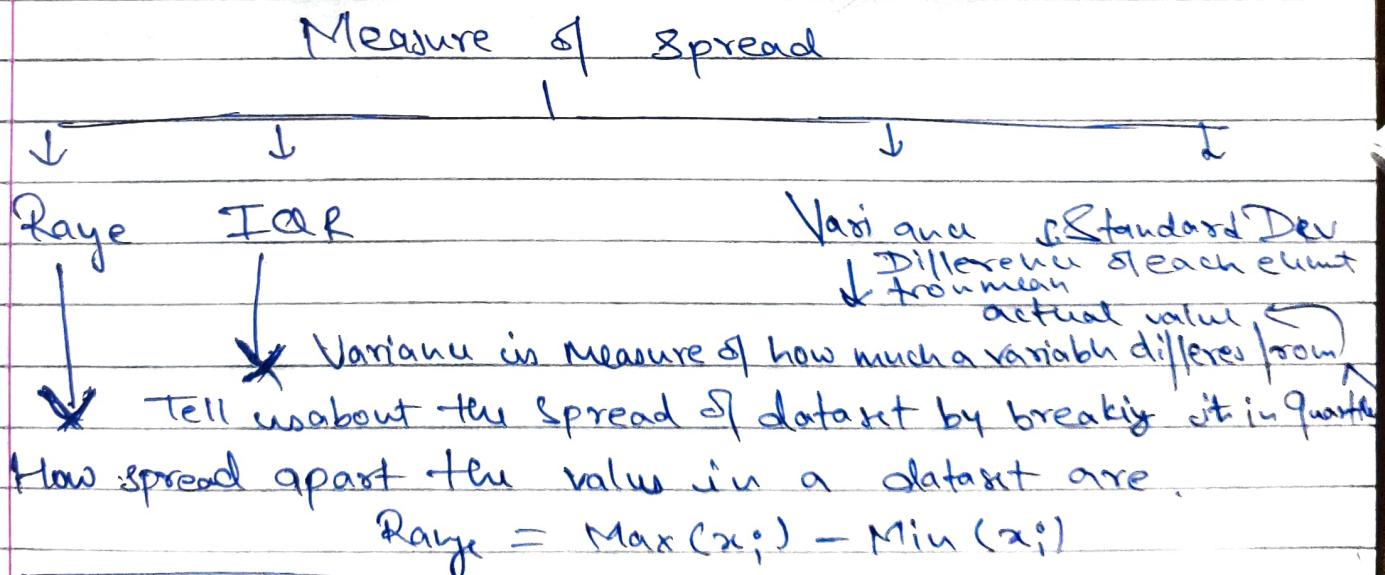
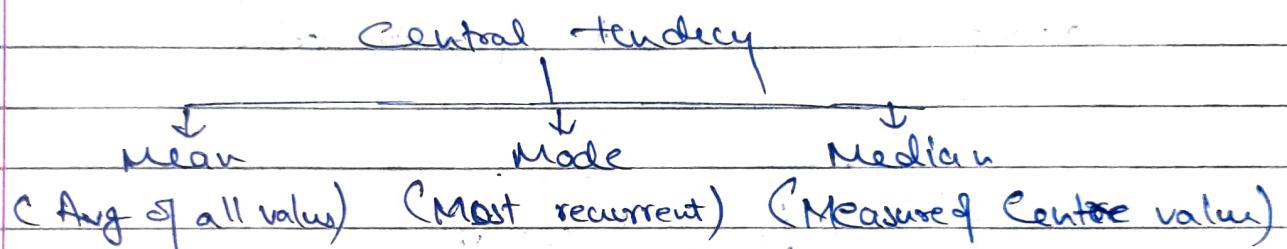
- used to use the data to provide description of the population, either through numerical calculation or graph or table.
- focuses on main characteristics of the data.

### ② Inferential Statistics

- It makes inferences & predictions about a population, based on a sample of data taken from the population in question. It is not guaranteed.
- It generalizes the large data set & applies probability to draw conclusion.

## Descriptive Statistics

- It is a method to describe & understand the features of a specific dataset by giving short summaries about the sample & measure the data.
- Descriptive statistics are broken down into two categories
  - Measure of central tendency (rep. Summary of Data set)
  - Measure of variability (spread)



Quartiles

$Q_1 \Rightarrow 25\text{th}$  value

$Q_2 \Rightarrow 50\text{th}$  value

$Q_3 \Rightarrow 75\text{th}$  value

$IQR = Q_3 - Q_1$

Standard Deviation =  $(x_i - \bar{x})^2$

$\bar{x}$  = individual data points

n = No of data points

$\bar{x}$  = mean.

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

Kash

Page No.:	
Date:	

$\mu = \text{Mean of Population}$   $N = \text{no of DP in Population}$

$\bar{x} = \text{Mean of Sample}$   $n = \text{no of DP in Sample}$

$\Rightarrow$  Population Variance is avg of Squared Deviation

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

$\Rightarrow$  Sample Variance is also avg of squared differences from the mean.

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

Standard Deviation is measure of dispersion of a set of data from its mean.

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

## Information Gain & Entropy.

Entropy : Measure of uncertainty in data.

$$H(S) = - \sum_{i=1}^N p_i \log_2 p_i$$

$S$  = Set of dataset all instances in the dataset.

$N$  = Number of distinct class values

$p_i$  = Probability event. different class value

Information Gain : How much information a particular feature / variable gives us about final outcome.

$$\text{Gain}(A, S) = H(S) - \sum_{j=1}^{|S_j|} \frac{|S_j|}{|S|} \cdot H(S_j)$$

$H(S)$  - Entropy of whole datasets

$|S_j|$  - No. of instances with j value of a attribute A

$|S|$  - Total no. of instances in dataset A.

$\vee$  - set of distinct value of an attribute A.

$H(S_j)$  - Entropy of subset of instance for an attribute A.

$H(A, S)$  - Entropy of an attribute A.

## Use case - Decision Tree.

$\Rightarrow$  Confusion Matrix

- It is a matrix used to describe the performance of a classification model or classifier. + one set of test data for which the true values are known + various additional data for testing.
- Represents a tabular representation of Actual vs Predicted values

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

True Positive (TP) + Predicted = 100%

Actual	No	Yes	Total
Predicted	TP	FP	TP + FP
No	FN	TN	FN + TN
Yes	FP	TP	FP + TP

Predicted	0	1
0	TP	FN
1	FP	TN

## Probability

$$P(x) = \frac{\text{Desired outcome}}{\text{Total}}$$

### Terminologies:

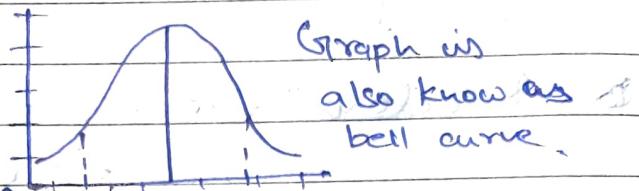
- **Random Experiment:** Experiment for which the outcome cannot be predicted with certainty.
  - **Sample Space:** The entire possible set of outcomes of a random experiment is called as the sample space ( $S$ ) of that experiment.
  - **Event:** Outcomes of an experiment
- Event Type:**
- Disjoint: Don't have a common outcome.
  - Non-Disjoint: Can have common outcome.

## Probability Distribution

### ① Probability Density Function.

→ It is the equation describing a continuous probability distribution.

Properties:

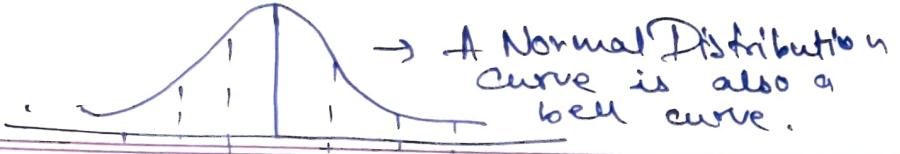


- ① Graph of PDF is continuous in range
- ② Area bounded by the curve of density function & the x axis is equal to 1 (area below the curve)
- ③ Probability that a random variable assumes a value between a & b is equal to the area under (PDF) bounded by a & b.

### ② Normal Distribution (Gaussian Distribution)

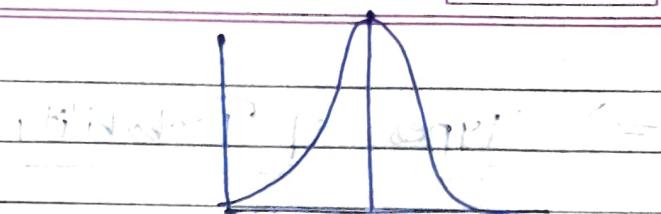
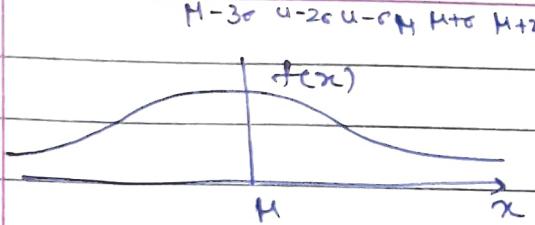
- Normal Distribution is a probability distribution that associates the <sup>normal</sup> random variable  $x$  with cumulative properties
- It denotes the symmetric properties of the mean. Idea behind is that the data near the mean occurs more frequently than the data away from the mean, or data near the mean represents the entire dataset
- Graph of Normal Distribution depends on two factors: the Mean & the Standard Deviation.

Mean: determines location of center of the graph  
 Standard deviation: determines height of the graph.



Page No.:  
Date:

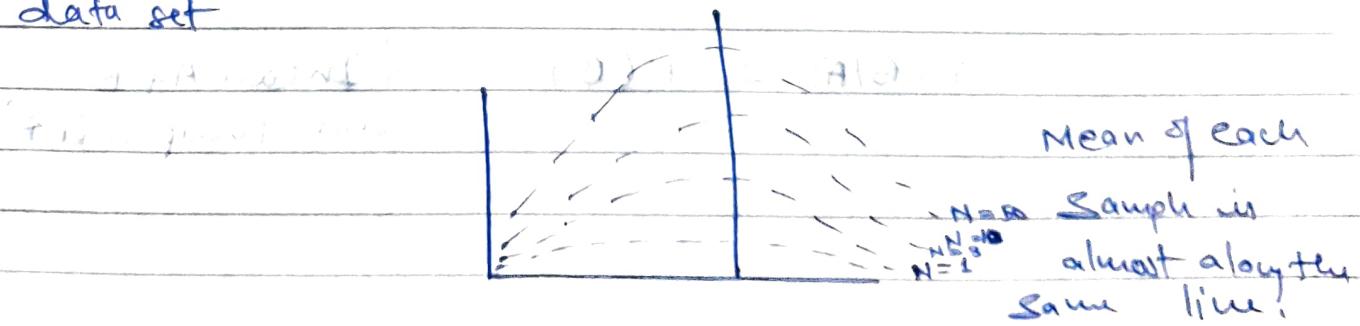
KOBY



If Standard Deviation is large, curve is short & wide or if standard deviation is small, curve is tall & narrow

⇒ Central Limit Theorem (only true for large dataset as small dataset has a lot of variations due to sample size factors)

- It states that the sampling distribution of the mean of any independent random variable will be normal or nearly normal if the sample size is large enough.
- If we have a large population & divided it into many samples, then the mean of all the samples of the population will be almost equal to the mean of the entire population. Hence, mean is normally distributed. Holds true only for a large data set



The accuracy or resemblance depends on two factors:

- ① No. of sample points considered
- ② Shape of underlying population  
→ shape depends on

- ① Standard Deviation
- ② Mean of sample.

## ⇒ Types of Probability

① Marginal Probability.

→ Probability of occurrence of single event.

$$\text{No. of hearts in a deck of card} = \frac{13}{52}$$

② Joint Probability.

→ Measure of two events happening at same time.

$$\text{Probability of getting Ace of hearts} = \frac{1}{52}$$

③ Conditional Probability, probability of an event given that another event has already occurred.

→ Probability of an event based on occurrence of a previous event.

An event B will occur given that A has already occurred.

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

(Events A & B are independent)

## ⇒ Baye's Theorem

- Shows the relationship between one conditional probability & its inverse.
- It is the probability of an event occurring based on prior knowledge of conditions that might be related to the same event.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

P(A|B) is referred to as posterior probability which means the probability of occurrence of A given B.  
 P(B|A) is referred to as likelihood ratio which measures the probability of occurrence of B given A.  
 P(A), is referred as Prior probability which represent the actual probability distribution of A.

Cq: we have 3 bowls.

Bowl A = 2 blue + 4 red

Bowl B = 8 blue + 4 red

Bowl C = 1 blue + 3 red.

1 ball is drawn from each. Probability of drawing a blue from A given 2 blue is drawn from B & C, are drawn in total.

Sol

Event 1 = Pickly a blue ball from bowl A.  $\Rightarrow A$

Event 2 = Pickly 2 blue balls.  $\Rightarrow X$ .

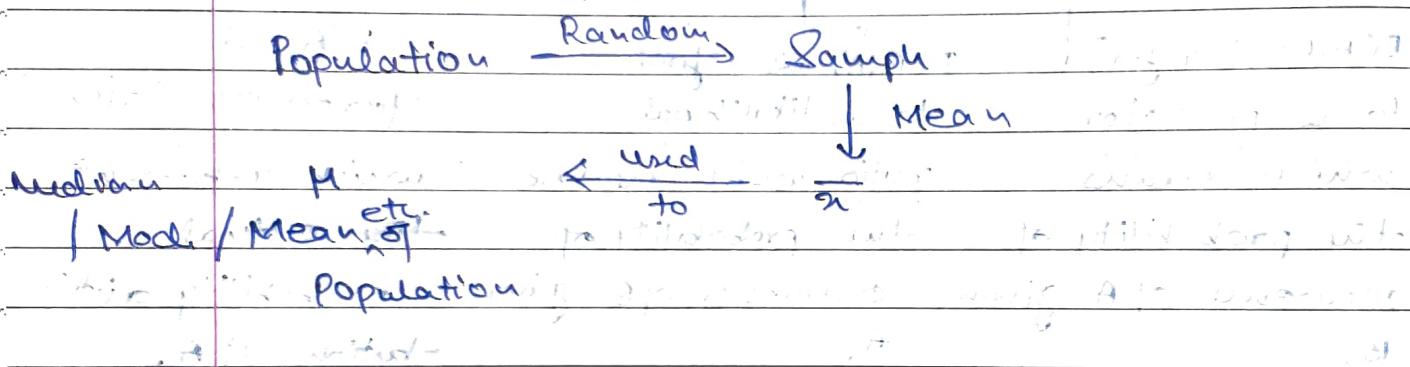
$$P(A/X) = \frac{P(A \cap X)}{P(X)}$$

$$P(A \cap X) = \frac{{}^6C_2 + {}^{12}C_2}{{}^{18}C_2}$$

$$P(X) = \frac{{}^6C_1 \times {}^{12}C_1 + {}^{12}C_1 \times {}^4C_1 + {}^4C_1 \times {}^6C_1}{{}^{22}C_2}$$

## $\Rightarrow$ Inferential Statistics & Point Estimation

Point Estimation is concerned with the use of the sample data to measure the single value that serves as an approximate value or the best estimate of an unknown population parameter.



Estimator = function of the sample used to find out the estimate.

$$\text{Ex: } \bar{x} \text{ Sample mean}$$

Estimate = pre-calculated value of result of the function

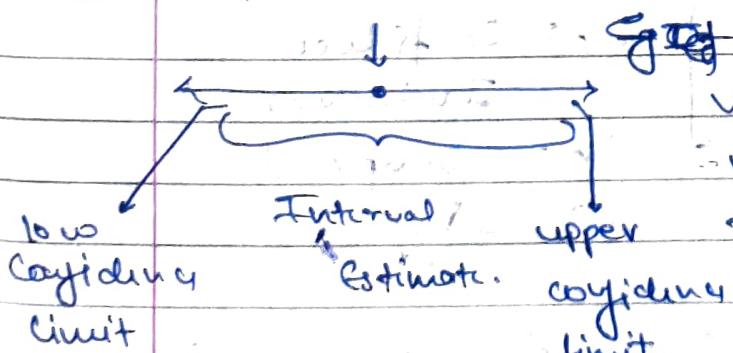
$$\text{Ex: population mean.}$$

## Method to find Estimator.

- ① **Method of Moments:** form an equation in the sample statistic by comparing it with the population parameter. Then solve the equations to get the estimates.
- ② To find out we take the known facts about the population. The population parameters are extended to the sample of size n. Then predict the known parameters such as mean, mode etc.
- ③ Maximum of Likelihood: estimating as below
- ④ Baye's Estimator
- ⑤ Best Unbiased Estimators.
- ⑥ Interval Estimation Method (Imp)

**Interval Estimate:** An ~~estimate~~ interval; a range of values used to estimate a population parameter with a degree of confidence.

Point Estimate



**Ex:** If we want to find a value (Mean, mode etc) of a population, we set a range & at that range the value lies.

So, instead of predicting a point we predict a ~~point~~ range, in which a value may occur.

- Difference b/w point estimate & actual population parameter value is Sampling Error
- When  $\mu$  is estimated, Sampling error =  $\bar{x} - \mu$

### Confidence Interval

- It is the measure of ~~your confidence about it~~ that the interval estimated contains the population parameter.
- It is used to describe the amount of uncertainty associated with a sample estimate of a population parameter.
- A range of values so constructed that there is specified probability of including the true value of a parameter within it to various levels of certainty.

### Margin of error - E

- For a given level of confidence it is the greatest possible distance b/w the point estimate & the value of parameter it is estimating.

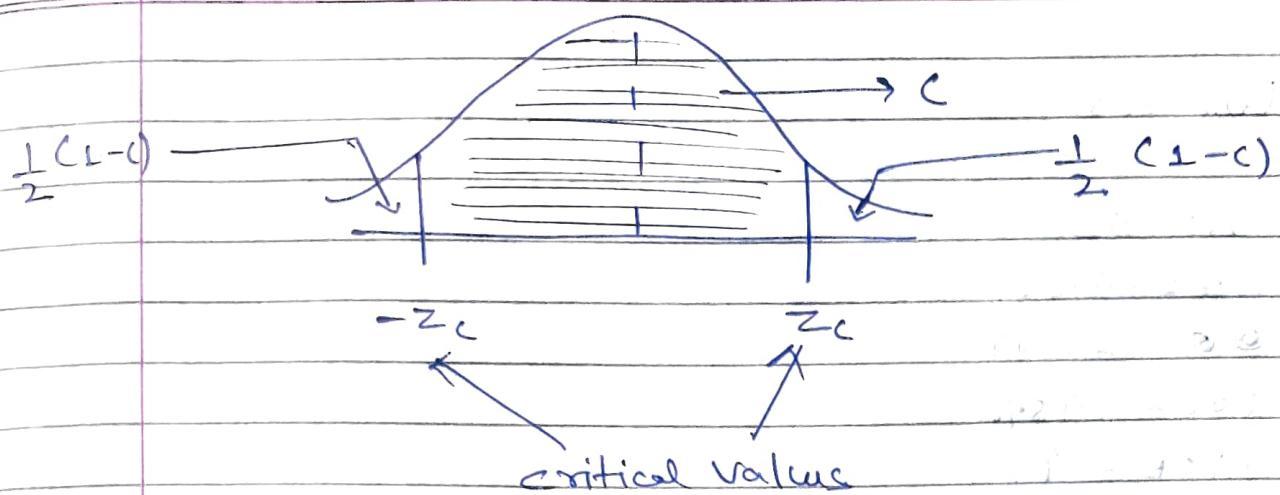
$$\text{Margin of error } E = Z_c \frac{\sigma}{\sqrt{n}}$$

$Z_c = \text{Confidence Interval}$

$n = \text{Sample size}$

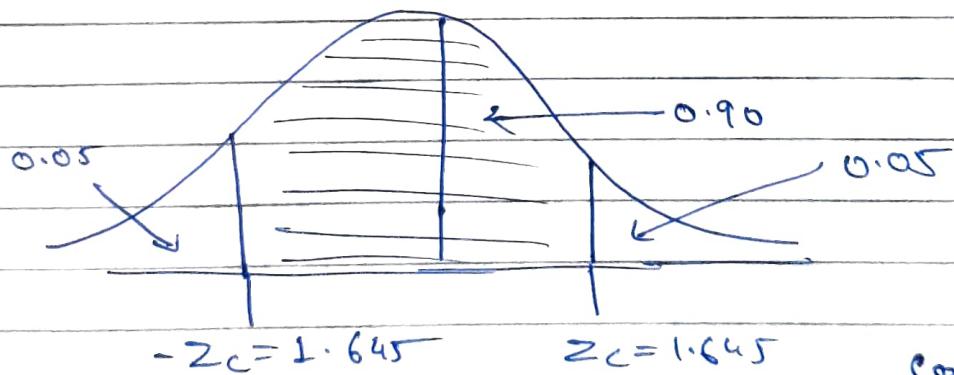
$\sigma = \text{SD}$

The level of confidence C is the probability that the interval estimate contains the population parameter.



c is the area beneath the normal curve b/w the critical values. Z-score can be calculated using standard normal table.

4 level of confidence = 90%. i.e. 90% confident that interval contains population parameter.



correspondingly  
Z score =  $\pm 1.645$

### Sequence :

- ① Sample statistic (used to calculate population parameter)
- ② Select confidence level.
- ③ Margin of error
- ④ Specify confidence interval

Use Case

$$n = 32$$

$$\bar{x} = 74.22$$

$$\sigma = 23.44$$

$$(Z) L.O.C = 957.$$

$$(E) M.O.E = 9$$

$$Q \Rightarrow f = 2, \frac{\sigma}{\sqrt{n}}$$

$$f = 1.96 \frac{23.44}{\sqrt{32}} \approx 8.12$$

## Hypothesis Testing

- Used to check whether hypothesis is accepted or rejected.

 Steps :

- ① State the hypothesis : Stating null & alternative hypothesis.
- ② Formulate an analysis plan: This stage involve construction of analysis plan.
- ③ Analyse Sample Data: Calculation & interpretation of the test statistic as described in the analysis plan.
- ④ Interpret Result: The application of the decision rule described in the analysis plan.

NULL HYPOTHESIS : Result is no different from ( $H_0$ ) assumption

ALTERNATIVE HYPOTHESIS : Result disproves the assumption.

Eg: 4 boys Nick, John, Bob, Harry are picked randomly each day. What is the probability that John is not cheating.

Sol

$$P(\text{John picked for a day}) = \frac{1}{4} = 0.25$$

$$P(\text{John picked for 3 days}) = \frac{1}{4} \times \frac{3}{4} \times \frac{3}{4} = 0.42$$

$$P(\text{John ————— 12 days}) = \left(\frac{1}{4}\right)^{12} = 0.032.$$

As,  $P(\text{Event}) < 5\% \text{ (threshold)}$ ; John is cheating.