

Descriptive Stats [Refers to the measure used to determine the centre of the distribution of data].

1. Measure of Central Tendency - {Mean, Median, Mode}
2. Measure of Dispersion - {Variance, S.D}   
↳ How they spread from the mean.

Statistics:- It is the science of collecting, Organising & Analyzing data.

↓  
Descriptive Stats - It consists of Organising & Summarising data

Interpretive Stats - Technique where in we use the data that we have measured to form conclusion.

### Sampling Techniques

#### 1. Simple Random Sampling



### Variance

Population

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

Sample

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

1. Simple Random Sampling :- Randomly they selecting of data points from sample.
2. Stratified Sampling - Where the population (N) is split into non-overlapping groups (strata). [Individual data].
3. Systematic Sampling - Population (N) -  $\frac{N}{k}$ <sup>th</sup> individual
4. Convenience Sampling - Availability, willingness or ease of access, and then they will select samples.

Variable :-

- Quantitative → Measured by numerical
- Qualitative → Categorical Variable

Measure of Variable Types:-

1. Nominal - Categorical data
2. Ordinal - Categorical data (Nominal) data but it follows Order
3. Interval -  $1-5, 5-15, 16-30$
4. Ratio - gap between Intervals Equal  
 $1-5, 6-10, 11-15$

phone  
↑

Decrease  
Continuous  
↓

Intervals or  
Discrete.  
↓

Only Decrease  
(Decrease Contn.)

## Percentiles & Quartiles

Percentage:- 1, 2, 3, 4, 5

% of the numbers that are odd

$$\% = \frac{\text{No. of numbers that are odd}}{\text{Total Number}} \times 100$$

$$= \frac{3}{5} \times 0.6 = 60\%$$

Percentiles (Crack, CAT, GMAT, SAT)

↳ A percentile is a value below which a certain percentage of observation lie.

Dataset- 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

What is the percentile ranking of 10?

$$\text{percentile Rank of } 10 = \frac{\text{No. of Values below } 10}{N} \times 100$$

$$\begin{aligned} & \frac{16}{20} \times 100 \\ & = 80\% \end{aligned}$$

What value exists at percentile ranking of 25%?

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

$$= \frac{25}{100} \times (20+1) = \frac{25}{100} \times 21 = 5.25$$

↓  
⑤ → 25%

### Five Number Summary

1. Minimum

2. First Quartile ( $Q_1$ ) (25%)

3. Median - (50%)

4. Third Quartile ( $Q_3$ ) (75%)

5. Maximum

IQR ]

\* Spread of the middle 50% of a dataset.

\* Is a key statistic in detecting Variability and Outliers.

### Removing the Outliers

{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27}

$$Q_1 = 3 \quad Q_3 = 7$$

$$IQR = Q_3 - Q_1 = 7 - 3 = 4$$

[ Lower fence  $\rightarrow$  Higher fence ]

$$[Q_1 - 1.5(IQR)] \quad [Q_3 + 1.5(IQR)]$$

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

$$= 3 - 1.5(4) = 3 - 6 = -3$$

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

$$= 7 + 1.5(4) = 7 + 6 = 13$$

$$[-3, 13]$$

## Removing data

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 22

↓  
Removed

Minimum = 1

Q<sub>1</sub> = 3

Median = 5

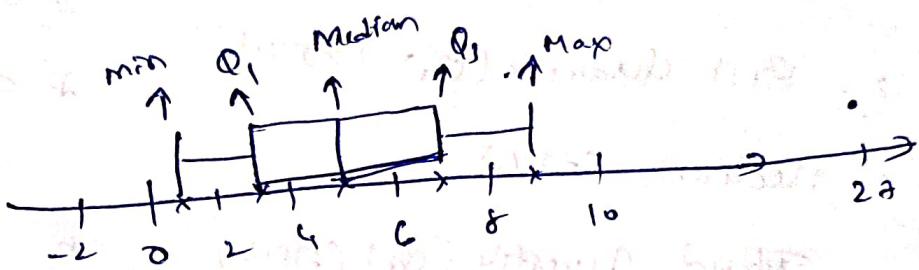
Q<sub>3</sub> = 7

Max = 9

→ 5 Number Summary



Box plot



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

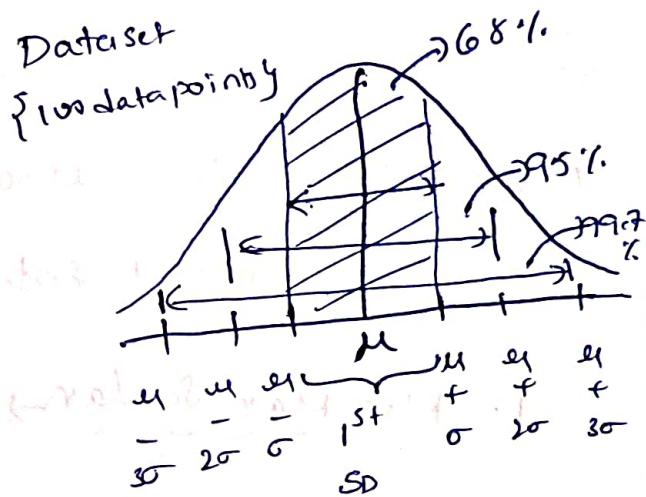
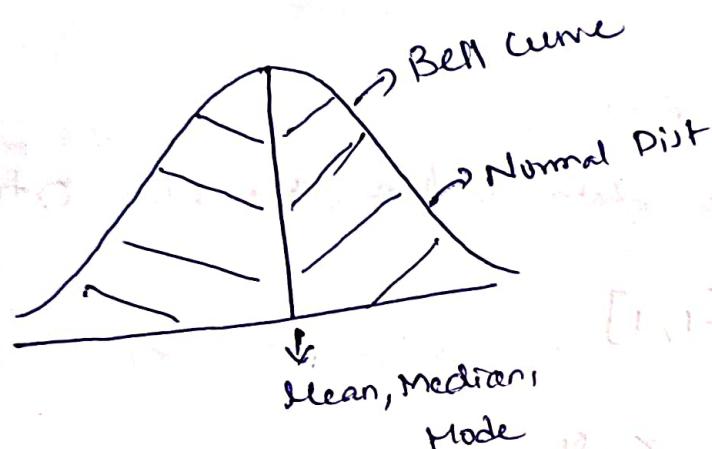
↳ Degrees of freedom

To correct bias and get an unbiased estimate of the population Variance

## Distributions:-

- ↳ Normal Distribution
- ↳ Standard Normal Distribution
- ↳ Z Score
- ↳ Log Normal Distribution
- ↳ Bernoulli Distribution
- ↳ Binomial Distribution

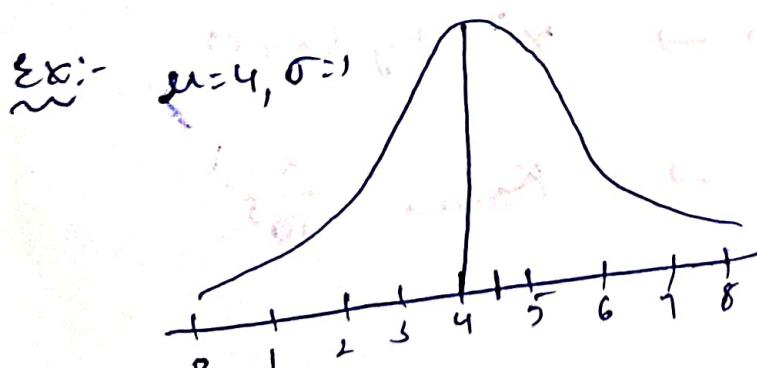
## ⑪ Gaussian / Normal Dist:-



### Empirical formula

68.1%, 95.4%, 99.7% Rule

↓  
This is present in 1st S.D re. Bell Shape



4.5 → Clear & knee

4.75 → Not clear

$$\begin{aligned} Z \text{ score} &= \frac{x - \mu}{\sigma} \\ &= \frac{4.75 - 4}{1} \\ &\approx 0.75 \text{ std} \end{aligned}$$

$$\mu=4, \sigma=1$$

$$\hookrightarrow \{1, 2, 3, 4, 5, 6, 7\}$$

$$[Z \text{ score} = \frac{x_i - \mu}{\sigma}]$$

$\hookrightarrow$  Normal Distribution / Gaussian Distri

$$\downarrow \{ -3, -2, -1, 0, 1, 2, 3 \}$$

$$Z(\mu=0) = \frac{1-4}{1} = -3$$

$\hookrightarrow$  Standard Normal Distribution

$$Z(2) = \frac{2-4}{1} = -2$$

### Standardization:-

The process of making data comparable by putting it on a common scale.

$$\mu=0, \sigma=1$$

### Normalization:-

$\hookrightarrow$  Convert entire data in ~~b/w~~ b/w 0 to 1

#### 1. Min Max Scalar $\rightarrow [-1, 1]$

$$2. Z \text{ Score} \rightarrow \frac{x - \mu}{\sigma}$$

$$3. \text{ Mean Normalization} \rightarrow \frac{x - \mu}{x_{\max} - x_{\min}}$$

$$4. \text{ Robust Scaling} \rightarrow \frac{x - \text{Median}(x)}{\text{IQR}(x)}$$

$$5. \text{ Log Transformation} \rightarrow x' = \log(x+c)$$

$$6. \text{ Decimal Scaling} \rightarrow x_{\text{scaled}} = \frac{x}{10^j}$$

## Probability :-

\* Is a measure of the likelihood of an event.

→ Additional Rule (probability, "or")

## \* Mutual Exclusive Event

→ Two Events are mutual Exclusive if they cannot occur at the same time.

Tossing a coin {H "or" T}  $\rightarrow p(H \text{ or } T) = p(H) + p(T)$

\* Also mutual Exclusive Event (probability, "AND") =  $1^2 / 2^2 = 1$

→ Multiple events can occur at the same time

Ex: Deck of cards { K, Q } ]

$$\therefore P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \text{ and } B)$$

$$\Rightarrow P(Q \text{ or } R) = P(Q) + P(R)$$

$$= \frac{9}{52} + \frac{13}{52} - \frac{1}{52} = \frac{16}{52}$$

- \* Naive Bayes (Conditional probability):
  - It is a probabilistic classification algorithm based on Bayes' Theorem.
  - It is widely used in Machine learning for tasks like Spam detection, Sentiment Analysis, text Classification.

Bayes' Theorem:

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

Where

$P(A|B)$  = probability of Event A happening given B has occurred (posterior probability)

$P(B|A)$  = probability of B happening given A has occurred (likelihood)

$P(A)$  = probability of Event A occurring (prior probability)

$P(B)$  = probability of Event B occurring (evidence)

Permutations:-

$$P(n,r) = {}^n P_r = \frac{n!}{(n-r)!}$$

A permutation is an arrangement of objects in a specific order.

Ex:- Passwords, seating arrangements, race rankings.

combinations:-

\* A combination is an selection of objects where the order does not matter.

$$C(n,r) = {}^n C_r = \frac{n!}{r!(n-r)!}$$

Ex:- Lottery, Selecting committees, teams.

\* Both are used for selecting of objects.

\* If order matters - Use permutations

\* If order doesn't matter - Use combinations.

## Hypothesis Testing:

- \* It is a statistical method used to make decisions about a population based on sample data.
- \* It helps us determine whether an assumption (hypothesis) about a dataset is true or not.

Step 1:-

Define the Hypothesis.

1. Null hypothesis ( $H_0$ ):-

The assumption that there is no significant effect

or no difference

2. Alternative Hypothesis ( $H_1$ ):-

The assumption that there is a significant

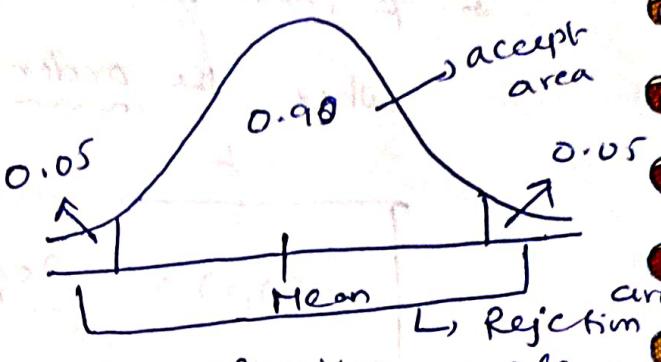
effect or difference

Step 2:- Select the Significance level ( $\alpha$ )

The significance level ( $\alpha$ ) is the probability of rejecting ( $H_0$ ) when it's actually true

\* Common choices: 0.05(5%) or 0.01(1%)

\* If  $\alpha = 0.05$ , we are 95% confident in our results



### Step-3: Choose Statistical Test

Type of data	Test to Use
1. Comparing Means	t-test (small datasets, samples)
2. Comparing two Groups	z-test (large samples)
3. Testing Relationships	Chi-Square test ( $\chi^2$ ) Correlation (or) Regressions.

### Step-4: Compute the Test Statistic & P-Value

Test Statistics :- A value that helps determine if the sample data is different from the assumption in  $H_0$ .

P-Value :-

The probability of observing the data if  $H_0$  were true.

### Step-5: Make a Decision [Reject or Accept (fail to Reject) $H_0$ ]

- \* If  $P\text{-Value} \leq \alpha \rightarrow \text{Reject } H_0$  [There is enough evidence to ~~fail~~ support  $H_0$ ]
- \* If  $P\text{-Value} > \alpha \rightarrow \text{fail to reject } H_0$  [there is not enough evidence to support  $H_1$ ].

Type I & II Errors

Rejecting a true Null hypothesis

		$H_0$ is true	$H_0$ is false
Decision			
Reject $H_0$	Reject $H_0$	Type I Error (False Positive) X	✓ Correct Decision (True Positive)
	Fail to Reject $H_0$	✓ Correct Decision (True Negative)	Type II Error (False Negative) X

Acknowledging a false null hypothesis.

predicted	Actually positive	Actually negative
	positive	negative
positive	True positive (TP)	false positive (FP) Type-I Error
negative	false Negative (FN)	True Negative (TN) Type-II Error

## One-Tailed Test (Directional Test)

- \* One-tailed test is used when we are testing for a specific direction (either an increase or decrease).
- \* We check if the sample mean is significantly greater than (or) less than the hypothesized mean, but not both.

Ex:-

- \* A railway company claims the average delay of trains is less than 10 min after introducing a new system. We need to do a hypothesis test.
- Hypothesis:-
- Null Hypothesis ( $H_0$ ) :- The average delay is 10 min or more.
- Alternative Hypothesis ( $H_1$ ) :- The average delay is less than 10 min.

When to use:-

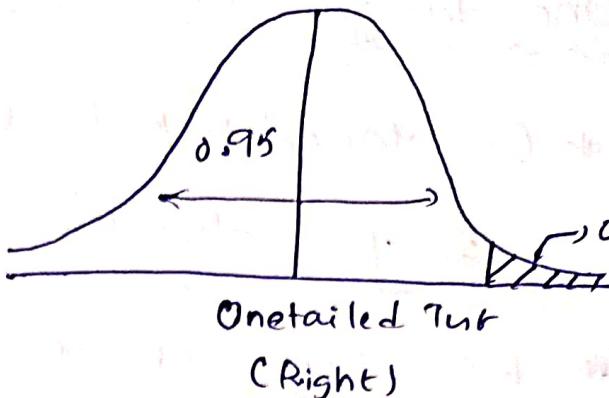
- \* If we expect only an increase or only a decrease (not both).

## Types of One-tailed Test

- \* Right tailed test
- \* Left tailed test.

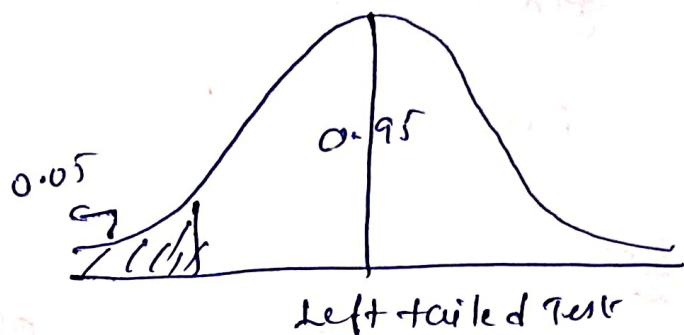
## \* Right tailed Test:-

→ Testing if the Sample mean is greater than the hypothesized mean.



## \* Left tailed Test:-

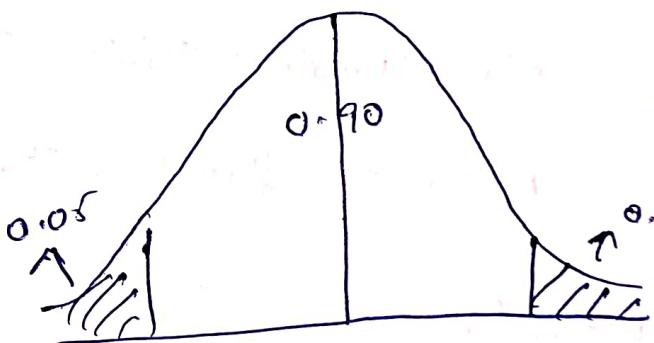
→ Testing if the sample mean is less than the hypothesized mean.



## Two tailed Test (Non-Directional Test)

\* A two-tailed Test is used when we are testing

for any significant difference  
(Either an increase or  
decrease).



\* we check if the sample mean is significantly different from the hypothesized mean, without specifying direction.

Example:-

- \* A ~~new~~ railway claims the average delay is 10 min after introducing a new system.

Hypothesis:-

- \* Null Hypothesis ( $H_0$ ) :- The average delay is 10 min
- \* Alternative hypothesis ( $H_1$ ) :- The average delay is not 10 min (could be more or less)

when to use:-

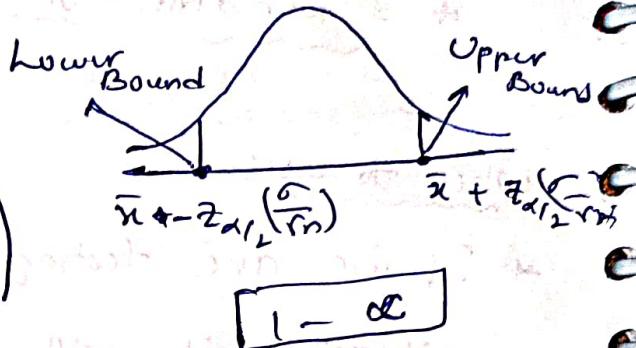
- \* If we are testing for any change, without specifying whether it will be higher or lower.

## Confidence Interval

- \* It is a range of values that likely contains the true population parameter (Eg., mean or proportion) with a certain level of confidence.
- \* It helps in Estimating the uncertainty of a Sample Statistic.

Formula :-

$$CI = \bar{x} \pm z_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right)$$



Here,  $\bar{x}$  = Sample mean

$z$  =  $z$  score (based on confidence level)

$\sigma$  = population standard deviation

$n$  = sample size

Example:-

1. Railway Monitoring - Estimating average train delay with confidence intervals

2. Manufacturing - Checking if the average weight of a ~~population~~ product is within an acceptable range

## Point Estimate

The value of any Statistic that Estimates the Value of a parameter.

Ex:- Inferential Stats

Sample mean  $\rightarrow$  population mean

$$\bar{x} \rightarrow \mu$$

$$\therefore \bar{x} = 2.9 \quad \mu = 3$$

Margin of Error.

## Confidence Intervals

point estimate  $\pm$  margin of error

$$\bar{x} \pm z_{\alpha/2} \left[ \frac{\sigma}{\sqrt{n}} \right]$$

Standard Error

## T-test :-

\* Used to compare the means of two groups when the sample size is small ( $n > 30$ ) and population variance is unknown.

### Types of T-test :-

→ One Sample T-test :- ( $\bar{x}$ )

\* Compare the mean of a sample to a known population mean.

Formula :-

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

Ex :- The company has a sample of 30 employees and wants to check if their average productivity score ( $\bar{x}$ ) differs from the industry standard ( $\mu$ ).

### Hypothesis :-

\*  $H_0$  : There is no significant difference ( $\mu = 70$ )

\*  $H_1$  : ~~The~~ The employee's productivity is different from 70 ( $\mu \neq 70$ )

where,  $\bar{x} = 75$  (Sample mean)

$\mu = 70$  (Population mean)

$s$  = Sample Standard deviation

$n = 30$  (Sample size).

→ Two Sample (Independent) Test:  
 & Compare the means of two independent Groups.

formula:-

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Ex:- The Company wants to compare IT employees & HR Employees, to see their productivity levels are different.

Hypothesis:-

- \*  $H_0$ : There is no significant difference ( $\mu_1 = \mu_2$ )
- \*  $H_1$ : Productivity levels are different ( $\mu_1 \neq \mu_2$ )

Here,  $\bar{X}_1, \bar{X}_2$  = Sample means of both groups.

$s_1^2, s_2^2$  = Sample Variances

$n_1, n_2$  = Sample sizes.

## Paired T-test (Dependent)

\* Compare the means of the same group before & after treatment.

formula

$$t = \frac{\bar{D}}{S_D / \sqrt{n}}$$

Ex:- The company tests if employee's productivity improved after a training program.

Hypothesis:-

\*  $H_0$  : No improvement after training ( $\mu_{\text{before}} = \mu_{\text{after}}$ )

\*  $H_1$  : productivity improved ( $\mu_{\text{before}} \neq \mu_{\text{after}}$ )

Here,

$\bar{D}$  = Mean of the differences (Before - After)

$S_D$  = Standard deviation of differences

$n$  = Number of pairs.

Scenario	(n < 30)	Test to use
Comparing one group to known population mean (population variance is unknown)		One-Sample T-test
Comparing two independent groups		Two Sample <del>T</del> -test (Independent Test)
Comparing same group before & after treatment		Paired T-test

## Z-Test :-

- \* Used to compare the means of two groups when the sample size is large ( $n > 30$ ) and the population variance is known.
- \* Used to test population proportion.

## Types of Z-Test:-

→ One Sample Z-Test:-  
 Compare the mean of a sample to a known population mean when the population variance ( $\sigma^2$ ) is known and sample size is large ( $n > 30$ ).  
 formula,  $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$

## Two Sample Z-Test:-

Compare the means of two independent groups when population variance is known & sample size is large.

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

Scenario ( $n > 30$ )	Test to use
Comparing one group to known population mean (Population Variance is known)	One Sample Z-test
Comparing two independent groups	Two Sample Z-test

## Anova Test:-

- \* Used to compare means of three or more groups to check if at least one group is significantly different.

## Types of Anova:

### One way Anova:

- \* Compares the means of more than two independent groups to determine if at least one group is significantly different.

Formula,  $F = \frac{\text{Between-Group Variance}}{\text{Within-Group Variance}}$

- \* Between Group Variance - Measures how much the group means differ from the overall mean.

$$MSB = \frac{\sum n_i (\bar{x}_i - \bar{x})^2}{k-1}$$

Measures how much the individual data points differ from their respective group mean.

$$MSW = \frac{\sum (x_{ij} - \bar{x}_j)^2}{N-k}$$

If F-statistic is greater than the critical value from the F-table we reject H<sub>0</sub>.

## → Two-way Anova (without Interaction):

\* Compares the means of groups across two independent Categorical Variables, without considering their Interaction.

Formula,  $F = \frac{\text{Variance Due to Factor A or B}}{\text{Error Variance}}$

Ex:- The company now wants to check two factors.

- 1. Work Environment (Remote, Hybrid, Office)
- 2. Experience level (Junior, Mid-level, Senior)

→ They want to analyze whether productivity is affected by work environment, experience level or both.

## → Two way Anova (with Interaction Effect):

\* In addition to analyzing the individual effect of two factors, it also checks if the two factors interact with each other.

formula,  $F = \frac{\text{Variance Due to Interaction}}{\text{Error Variance}}$

Ex:- The company now wants to check whether experience level influences the effect of work environment on productivity.

→ They want to analyze

- 1. main effect of work environment

- 2. Main effect of Experience level

- 3. Interaction effect b/w Work Environment & Experience level

Scenarios	Test to Use
comparing one categorical factor with more than two groups	One-way Anova
comparing two categorical factors independently	Two-way Anova (without interaction)
checking if two categorical factors interact with each other	Two-way Anova (with interaction)

## Chi-Square Test :-

The Chi-Square ( $\chi^2$ ) test is used for Categorical data to check whether there is a Significant association b/w two variables or whether the observed data follows an Expected distribution.

### 1. Chi-Square Goodness of fit test:-

To check if the observed frequency distribution of a categorical variable matches an expected distribution.

$$\chi^2 = \frac{\sum (O_i - E_i)^2}{E_i}$$

where,

$O_i$  = Observed frequency

$E_i$  = Expected frequency

### 2. Chi-Square Test for independence:-

To check if two categorical variables are independent or related

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} ; E_{ij} = \frac{\text{Row total} \times \text{Column total}}{\text{Grand total}}$$

#### Scenario

Checking if observed data follows an expected distribution

Checking if ~~two~~ two categorical variables are related

#### Test to use

Chi-square Goodness of fit

Chi-Square Test for independence

## Scenario

## Test to Use

Comparing means of two groups, Small Sample size

T-test

Comparing means of two groups, Large sample size

Z-test

Comparing means of 3+ groups

ANOVA test

Testing relationship b/w two categorical Variables

Chi-Square test

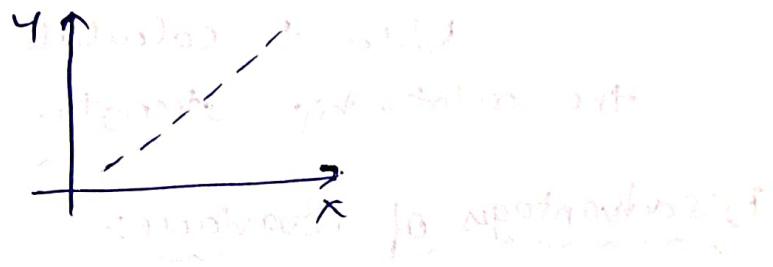
## Covariance:-

\* Covariance is a statistical measure that indicates the direction of the linear relationship between two variables.

## Types:-

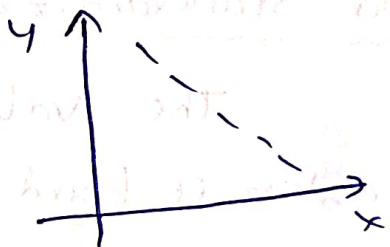
### \* Positive Covariance:-

When one Variable Increases, the other also tends to Increase.



### \* Negative Covariance

When one Variable increases, the other tends to decrease.



### \* Zero Covariance

No relationship between the Variables.



## Formula,

for population,

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

$\text{cov}(x,y) =$

for Sample,

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

## Advantages of covariance:-

### \* Indicates Relationship Direction:-

Shows whether two variables move together (positive or negative relationship)

### \* Useful in Portfolio Management:-

Helps in diversifying investments by analyzing how asset prices move together.

### \* Foundation for Correlation:-

Used to calculate correlation, which normalizes the relationship strength.

## Disadvantages of covariance:-

### \* No Standardized Scale:-

The value depends on the units of the variables, making it hard to interpret.

### \* Doesn't Show Strength:-

A large covariance doesn't always mean a strong relationship.

### \* Sensitive to Outliers:-

Extreme values can distort the covariance calculation.

## Correlation:-

- \* Correlation measures both the direction & strength of the relationship between two variables.
- \* It is a standardized version of covariance, making it easier to interpret.

Formula :-

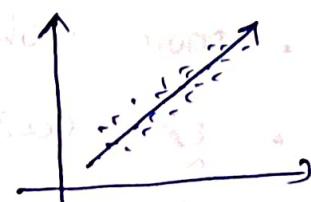
$$\gamma = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

- \* It is derived from covariance and ranges between  $-1$  &  $1$ .

- \* Unlike covariance, which only indicates the direction of the relationship, correlation provides a standardized measure.

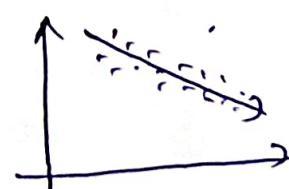
- Positive correlation (close to +1) :-

As one variable increases, the other variable also increase.



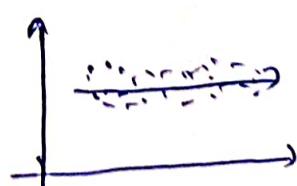
- Negative correlation (close to -1) :-

As one variable increases, the other variable also decrease.



- Zero correlation :-

There is no linear relationship between the variables.



## Types of Correlation:-

### \* Pearson Correlation (Linear Correlation):-

- Measure the linear relationship between two variables.
- Value ranges from -1 to +1.
- Ex:- If temperature increases, ice cream sales also increase.

### \* Spearman's Rank Correlation:-

- Measure the monotonic relationship (not necessarily linear).
- Works on ranked data instead of actual values.
- Ex:- If a student studies more, their exam rank improves.

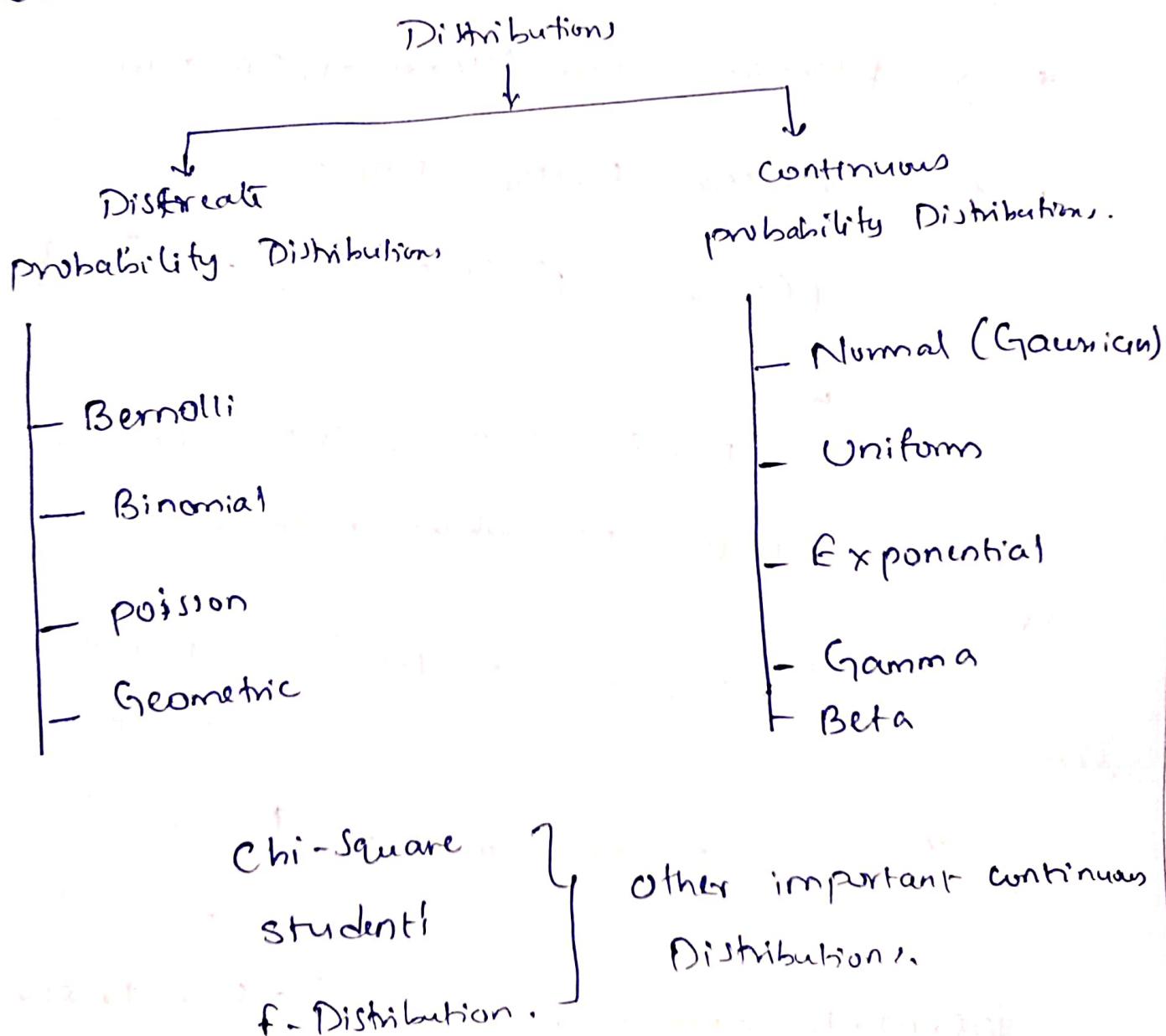
### \* Kendall's Tau Correlation:-

- Measure the association between two ranked variables.
- More robust for small datasets.
- Ex:- Customer Satisfaction ranking vs product sales ranking.

## Distributions:-

- It refers to how values of a dataset or a random variable are spread or arranged.
- It shows the possible values and their frequency of occurrence.
- It helps in understanding patterns, making predictions & applying statistical tests.

## Types of Distributions:-



# 1. Discrete probability Distributions

These are used for Variables that take  
Countable Values.

Ex:- Number of Success, number of Events.

a. Bernoulli Distribution:- (Single trial) [Success/Failure]

- \* Represents a Single trial with only two possible outcomes. Success (1) or failure (0)
- \* Ex:- flipping a coin (Heads = 1, Tails = 0)

\* Probability Mass Function (PMF):

$$P(X=x) = p^x (1-p)^{1-x}, \quad x \in \{0, 1\}$$

where,

$p$  - probability of success. ( $0 \leq p \leq 1$ ).

$1-p$  - probability of failure

if,

$$x=0, \quad P(X=0) = p^0 (1-p)^{1-0} = 1(1-p) \\ = 1-p.$$

$$x=1, \quad P(X=1) = p^1 (1-p)^{1-1} = p(1)$$

$$= p.$$

Bernoulli distribution,

$$P(X) = \begin{cases} 1-p, & \text{if } x=0 \\ p, & \text{if } x=1 \end{cases}$$

b. Binomial Distributions: (n-independent trials), count success

- The Binomial Distribution represents the probability of obtaining 'K' success in 'n' independent trials, where each trial is a Bernoulli trial with a success probability - P.

\* Probability Mass function (PMF)

$$P(X=K) = {}^n C_K P^K (1-P)^{n-K}, K=0,1,2 \dots n$$

where,

\*  ${}^n C_K = \frac{n!}{k!(n-k)!}$  is a Binomial Coefficient

which representing the number of ways to choose K success from n trials.

\*  $P =$  is the probability of success in Each trial

\*  $(1-P) =$  is the probability of failure in Each trial

Ex: flip a fair coin ( $P=0.5$ ) → 3 times ( $n=3$ )

calculate the probability of getting exactly 2 heads ( $K=2$ )

$$P(X=2) = \frac{3!}{2!(3-2)!} (0.5)^2 (1-0.5)^{3-2} = \frac{3!}{2!(1)!} (0.25)(0.5)$$

$$= \frac{3}{2} \times 0.125$$

Exactly 2 heads in 3 flips  $\Rightarrow 0.375$  (or) 37%  $\Rightarrow 0.375$

Position Distribution: (Time) Count Events in a time / Space Interval

- \* This models the probability of observing a certain number of events in a fixed interval of time or space.
- \* Given that the events occurs at a constant average rate & independently of each other

→ Probability mass function (PMF)

$$P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}, k=0, 1, 2, \dots$$

where,

$X \rightarrow$  is the number of occurrences of an event

$\lambda \rightarrow$  (Lambda) is the average number of occurrence per time / space unit (also called the rate parameter).

$e \rightarrow$  Euler's Number ( $\approx 2.718$ )

$k! \rightarrow$  factorial of  $k$ , possible arrangements of  $k$  events

Ex:- Call center receives an average of 5 calls per hour ( $\lambda=5$ ), and we want to find the probability of receiving exactly 3 calls in hour ( $k=3$ )

$$P(X=3) = \frac{5^3 e^{-5}}{3!} = \frac{125 \cdot e^{-5}}{6} = \frac{125 \times (2.718)^{-5}}{6}$$

∴ Exactly 3 calls in an hour is

$$\frac{125 \times 0.0067}{6} = 0.1396$$

d. Geometric :- (Number of trials until first success)

- The Number of Bernoulli trials needed to get the first success in a Sequence of independent trials,
- where Each trial has a Success probability  $P$  and a failure probability  $(1-P)$ .

\* Probability Mass function(PMF) :

$$P(X=k) = (1-P)^{k-1} \cdot P, \quad k=1, 2, 3, \dots$$

Where:

$x$  - is the number of trials until the first success occurs

$p$  - probability of success in a single trial.

$1-p$  - probability of failure

$k$  - The trial number where the first success happens

Ex:- flip a biased coin, where the probability of getting heads (success) is  $P=0.4$ . we want to find the probability of that the first heads appear on the 3rd flip.

$$\begin{aligned} P(X=3) &= (1-0.4)^{3-1} (0.4) \\ &= (0.6)^2 \times 0.4 \\ &= 0.36 \times 0.4 \\ &= 0.144 \end{aligned}$$

∴ Probability of getting the first head on 3<sup>rd</sup> flip is  $0.1444$  ( $14.4\%$ )

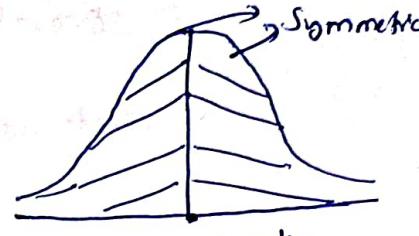
## Continuous Probability Distributions:-

These describe Variables that can take any value within a range (Ex:- Height, temperature, time).

### a. Normal Distribution (Gaussian) :- (Natural phenomena)

\* A bell-shaped curve where most values cluster around the mean.

Ex:- Heights of people, IQ scores.



Probability density function (PDF)

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, -\infty < x < \infty.$$

Normalizing factor  $\rightarrow$  Exponential part

where,  $x$  - Random variable (value we evaluate the distribution at)

$\mu$  = Mean

$\sigma$  = Standard deviation

$\sigma^2$  = Variance

e = Euler's no., ( $\approx 2.718$ )

$\pi = 3.1416$ .

### Properties:-

1. Symmetrical; Centered at  $\mu$  with equal probabilities on Both sides
2. Mean = Median = Mode,
3. 68-95-99.7 Rule (Empirical)

Ex:-  $\mu = 100, \sigma = 15, x = 120$

$$f(120) = \frac{1}{15(\sqrt{2\pi})} e^{-\frac{(120-100)^2}{2(15)^2}}$$

$$= \frac{1}{37.67} e^{-0.8889}$$

Probability density at  $x = 120$ , is  $\approx 0.0267$

## b. Uniform distribution:- (Equal Probability in a range)

\* The Uniform distribution is a probability distribution

where all outcomes are equally likely within a given range

\* It is divided into:

→ Discrete Uniform Distribution (finite set of outcomes)

→ Continuous Uniform Distribution (infinite set of outcomes with an interval)

→ Discrete Uniform Distribution:- (All probabilities are equal)

$$\text{PMF} \rightarrow P(X=x) = \frac{1}{n-1}, x = x_1, x_2, \dots, x_n$$

Ex:- Rolling a fair coin.

outcomes =  $\{x : \{1, 2, 3, 4, 5, 6\}\}$ , So,  $n=6$ .

$$P(X=3) = \frac{1}{6} = 0.1667 \text{ (for rolling a 3).}$$

→ Continuous Uniform Distribution:-

→ In which all values in a given range have equal probability.

$$\text{PDF} \rightarrow f(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$$

Here,  $a, b$  - minimum & maximum values

$\frac{1}{b-a}$  - Height of the distribution is constant

## Cumulative distribution function(CDF):-

$$F(x) = \begin{cases} 0, & x < a \\ \frac{x-a}{b-a}, & a \leq x \leq b \\ 1, & x > b. \end{cases}$$

Ex:-

Bus arrives randomly b/w 10 Am & 10:30 Am.

If you arrive at a random time in the interval,  
the waiting time ( $x$ ) is uniformly distributed b/w

$a = 0$  (minutes, 10:00 Am)  
 $b = 30$  (minutes, 10:30 Am)

PDF,  $f(x) = \frac{1}{30-0} = \frac{1}{30}$ ,  $0 \leq x \leq 30$

Find the Probability that wait less than 10 min.

CDF,  $P(x \leq 10) =$

$P(x \leq 10) = \frac{10-0}{30-0} = \frac{10}{30} = 0.3333$

Waiting less than 10 min = 33.33%.

c. Exponential Distribution - [Waiting times [time until the next event]]

→ Used to model the time until an event occurs in a Poisson process.

→ It is commonly used to represent waiting times such as the time b/w customer arrivals at a service station or the time until a machine fails.

PDF,  $f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$

Here, CDF,  $F(x) = P(X \leq x) = 1 - e^{-\lambda x}, x \geq 0.$

$x$  - time until the event occurs

$\lambda$  - rate parameter (inverse of mean,  $\lambda = \frac{1}{4}$ )

$e$  = Euler's number ( $\approx 2.718$ ).

Ex: Call center receives 5 calls per hour on average.

The time b/w calls follows an Exponential distribution

$\lambda = 5$  calls per hour

Find the probability that the next call occurs with

in ~~10 min~~  $(1/6 \text{ hour})$

$$\text{CDF} : P(X \leq \frac{1}{6}) = 1 - e^{-5(\frac{1}{6})}$$

$$= 1 - 0.4352$$

$$= 0.5648$$

$\therefore 56.48\%$  Probability that the next call arrives within 10 min

d. Gamma Distribution: [Generalization of Exponential]

\* Gamma Distribution used to model the time required for multiple independent events to occur in a point process.

\* It is a generalization of the exponential distribution which models the time until a single event occurs.

PDF;  $\rightarrow$

$$f(x) = \begin{cases} \frac{\lambda^k x^{k-1} e^{-\lambda x}}{\Gamma(k)}, & x \geq 0 \\ 0, & x \leq 0 \end{cases}$$

CDF,  $\rightarrow$   $F(x) = P(X \leq x) = \frac{\gamma(k, \lambda x)}{\Gamma(k)}$

$$\gamma(k, \lambda x) = \int_0^x t^{k-1} e^{-t} dt$$

$(\lambda + 1)^{-k}$

Ex:- Call center 5 call per hour on average,

three calls arrive follows gamma distribution.

$$K = 3, \lambda = 5$$

Find the Probability that three calls arrive within 30 min (1/2 hours)

$$P(X \leq 1/2) = \frac{\gamma(3, 5(1/2))}{\Gamma(3)}$$

$$= 0.875$$

87.5%, probability that third call arrives within 30 min

## Beta Distribution: [Probability distributions in Bayesian Stats]

→ It is used to model probabilities and proportions, often in Bayesian statistics & machine learning.

$$\text{PDF}, \quad f(x; \alpha, \beta) = \begin{cases} \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}, & 0 < x < 1 \\ 0, & \text{Otherwise} \end{cases}$$

$\alpha, \beta$  - Shape parameters (both must be greater than 0)

$B(\alpha, \beta)$  = Beta function, which normalizes the distribution

$$B(\alpha, \beta) = \int_0^1 t^{\alpha-1}(1-t)^{\beta-1} dt.$$

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}, \quad \Gamma(n) = (n-1)!$$

if n is integer

CDF,

$$F(x) = I_x(\alpha, \beta) = \frac{\int_0^x t^{\alpha-1}(1-t)^{\beta-1} dt}{B(\alpha, \beta)}$$