

29/10/2025

Introduction

- * It is used to check the quality of the data. (Descriptive Statistics)
- * To make a statement or a conclusion we use Inferential Statistics.

What is Statistics:-

It is the science of collecting, Organizing and Analyzing data. (For better decision making)

What is Data?

Data means facts or Piece of information that also can be measured.

Ex: The IQ of a class students

$$\{ 98, 97, 68, 57, 110 \}$$

We can measure Avg, Min, Max

Descriptive Statistics

It consists of organizing & summarizing Data.

Inferential Statistics

* Techniques where we used the Data that we have measured to form Conclusion.

* To make a statement / Conclusion on a Descriptive statistics we use inferential statistics.

1) Are the avg mark of the Java students is same as python class students in Besant?

→ Inferential Statistics

2) What is the average marks of SQL students?

→ Descriptive Statistics

Population (N) and Sample (n) :-

* The Entire group of the data we call it as Population (N).

Eg: All People in India

* A Subset of a population we call it as a Sample.

Eg: One lakh people from different region of India.

* Population are Larger than Sample.

* Sample should be representative of the population.

* Samples allow for easier, faster & less costly data collection.

Types of Sampling Techniques:-

1) Simple Random Sampling

Every member of a population has an equal chance of being selected for our sample.

Eg: - Avg mileage of any 120 cc bikes.

- Ratio of married people in banglore

2) Stratified Sampling:

Where the population is split into non overlapping groups.

Eg: i) The person is eligible for voting or not.

3) Systematic Sampling:

From the population, Every n^{th} sample we will collect.

Eg: i) While doing surveying in the mall on modernization, collecting information of every fifth person who is coming out from the mall.

4) Convenience Sampling

The sample is collected based on our convenience from the particular domain experts.



30/10/2025

Note:- Sampling techniques Selection always depends on Problem Statement.

Variable:

A Variable is a property that can take an any value.

Two kinds of Variables:-

i) Quantitative (Numerical) Variables

ii) Qualitative (Categorical) Variables

① Quantitative Variables:-

A Value can be measured and we can perform mathematical operation like (A,S,M,D)

Ex: mpg, weight, height ...

Continuous (float) $\rightarrow \infty$

Numerical

0-1

Discrete (int) $\rightarrow ^1$

Numerical: age (10, 20, 25)

Discrete

Categorical: gear (4, 5, 6)

② Qualitative / Categorical Variable:-

Non-Measurable Data and Based on Some characteristics we can define Categorical Variables.

Ex: gender

- male
- female
- other

DT

IT
non-IT

Blood

- A+
- B+
- A-
- B-
- O+
- O-

Variable Measurement Scales:-

4 Types of Measured Variables

1) Nominal Data

2) Ordinal Data

3) Interval Data

4) Ratio Data

1) Nominal Data:-

The categorical Data which are having different classes.

2) Ordinal Data:-

Order of the Data matters but values

Does not

Marks [50, 40, 100, 91, 35] {Rank is imp} 0 - 100

Marks [40, 80, 75, 60, 90] {Value is not} 0 - 100

3) Interval Data: Order matters and Value also matters but natural zero is not present.

4) Ratio Data:
The ratio data can measured, ordered, equi distant and have meaningful zero.

Eg: Sal, height, weight, ...

Descriptive Statistics:

I Measure of Central Tendency:

- 1) Mean
- 2) Median
- 3) Mode

1) Mean:

$$\text{Population Mean } (\mu) = \frac{\sum_{i=1}^N x_i}{N}$$

$$\text{Sample Mean } (\bar{x}) = \frac{\sum_{i=1}^n x_i}{n}$$

- Median
- Sort the value Either asc or des & order.
 - Choose the Mid value.
 - If you get 2 mid values take avg of those 2 values.

$[1, 2, 2, 3, 4, 5]$



Mean = 2.8

Median = 2.5

$[1, 2, 2, 3, 4, 5, 100]$



Mean = 16.7

Median = 3

100 = outlier

out of Boundary

- Mean will be affected by outliers whereas median won't affect by outliers.
- Used for Null value Imputation

3) Mode

Most Repetitive Value

$[1, 2, 3, 1, 2, 2, 3, 4, 5]$

Mode = 2



31/10/2025

	Person 1	Person 2
M	7 am	8
T	7.30 am	11
W	8 am	9
T	7.15 am	7
F	7.30 am	10
S	?	?

Manager ask Team lead at what time the persons will come?

$$\text{Person 1} \rightarrow 7-8 - 7.15(15\text{min}) \\ 7.30(30\text{min}) \\ 7.45(15\text{min})$$

$$\text{Person 2} \rightarrow 9-10 - 7(2\text{hr}) \\ 11(1\text{hr})$$

- If variance is high (8-11) Prediction accuracy low
- If variance is low (7-8) Prediction accuracy high.

(II)

Measure of Dispersion

1) Variance

2) Standard Deviation

3) Range

1) Variance:-

- Population Variance (σ^2)
- Sample Variance (s^2)

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

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

$\underline{+} + \underline{+} = 57 \rightarrow n=3$ $n-1 = \text{degree of freedom}$

$(20+30) + \underline{?} = 57$ \downarrow

$n-1$

$3-1=2$

Degree of freedom

$\{1, 2, 2, 3, 4, 5\}$ km

Population

Variance $\sigma^2 = \frac{\sum_{i=1}^6 (x(i) - \mu)^2}{6}$

$$\mu = \frac{1+2+2+3+4+5}{6} = 2.8$$

$$\sigma^2 = \frac{(1-2.8)^2 + (2-2.8)^2 + (3-2.8)^2 + (4-2.8)^2 + (5-2.8)^2}{6}$$

$$\sigma^2 = \frac{10.83}{6}$$

$$\sigma^2 = 1.8 \text{ km}^2$$

$$\sigma = 1.34 \text{ km}$$

2) Standard Deviation:-

$$SD = \sqrt{\text{Variance}}$$

Population

SD

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

$$\sigma = \sqrt{1.8}$$

$$\sigma = 1.34$$

Sample Variance (S^2)

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

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

$$S^2 = \frac{(1-2.8)^2 + (2-2.8)^2 + (3-2.8)^2 + (4-2.8)^2 + (5-2.8)^2}{5}$$

$$S^2 = 2.16 \text{ km}^2$$

Standard Deviation of Sample Variance

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

$$S = \sqrt{2.16}$$

$$S = 1.46$$

3) Range:-

Range = Maximum - Minimum

$$\text{Range} = 5 - 1 = 4$$

Percentile & Quantiles

Percentile \rightarrow It is a value below which a certain percent of observation will lie.

Eg: $\{1, 1, 2, 3, 4, 5, 5, 6, 7, 7, 8\}$

How much % of data will come below 6?

Percentile rank of $x = \frac{\# \text{ of value below } x}{N} \times 100$

$$\text{rank of } 6 = \frac{7}{11} \times 100$$

$$= 63\%$$

63% of observation data is ≤ 6

Quartile:-

It will help to find the value which is present in the given Percentile rank.

Eg: $\{1, 1, 2, 3, 4, 5, 5, 6, 7, 7, 8\}$

Which value is present at 25%?

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

$$= \frac{25}{100} \times 2$$

$$\boxed{\text{Value} = 3 \text{ in data}}$$

$$90\% = \frac{90}{100} \times 12$$

$$= 10.8 - \text{index}$$

$$= 10. (Index)$$

$$\therefore = 7$$



Five Number Summary :-

1) Minimum

2) First Quartile (Q_1) 25%.

3) Median (Q_2) 50%.

4) Third Quartile (Q_3) 75%.

5) Maximum

Note:- Choose these 5 numbers after removing the outliers from the data by finding boundary values.

[Lower fence upper fence]

$$LF = Q_1 - 1.5(IQR)$$

$$UF = Q_3 + 1.5(IQR)$$

IQR (Inter Quartile Range)

$$IQR = Q_3 - Q_1$$

$$\{ 1, 1, 2, 3, 4, 4, 5, 5, 6, 7, 7, 8, 8, 9, 28, 36 \} = 17$$

$$Q_1 - \text{Taille } \times (Q_1) = \frac{25}{100} \times (17+1)$$

$$Q_1 = 4.5 \rightarrow 4 \rightarrow 3$$

$$(Q_3) = \frac{75}{100} \times 18$$

$$Q_3 = 13.5 \rightarrow 13 \rightarrow 8$$

$$IQR = Q_3 - Q_1 = 8 - 3$$

$$\boxed{IQR = 5}$$

$$LF = 3 - 1.5(5)$$

$$\boxed{LF = -4.5}$$

$$UF = 8 + 1.5(5)$$

$$\boxed{UF = 15.5}$$

$$\therefore [-4.5 \text{ (Lower Fence)} \longleftrightarrow (Upper Fence) 15.5]$$

Anything out of the range of -4.5 & 15.5 is outlier.

\therefore In $\{ 1, 1, 2, 3, 4, 4, 5, 5, 6, 7, 7, 8, 8, 9, 28, 36 \}$

28 & 36 are outliers

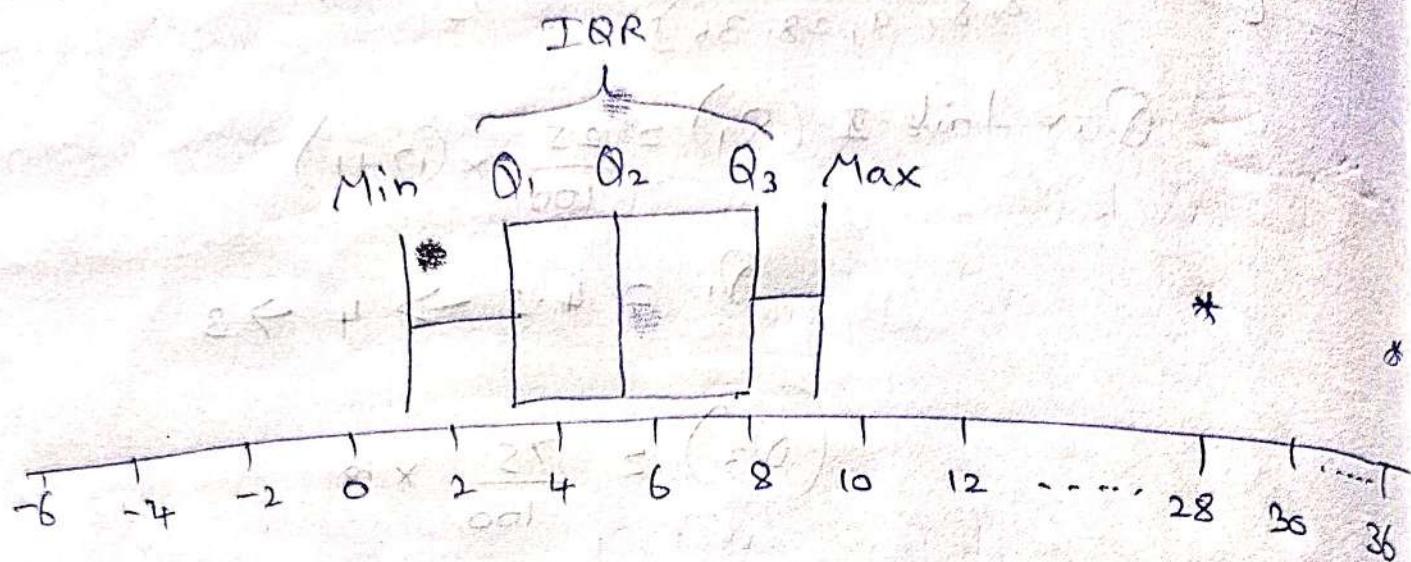
$$\text{Minimum} = 1$$

$$Q_1 = 3$$

$$\text{Median}(Q_2) = 5$$

$$Q_3 = 8$$

Box Plot or Five no Summary:-



This Graph is used to find the Outlier.

3/11/2025

Different types of Distribution:-

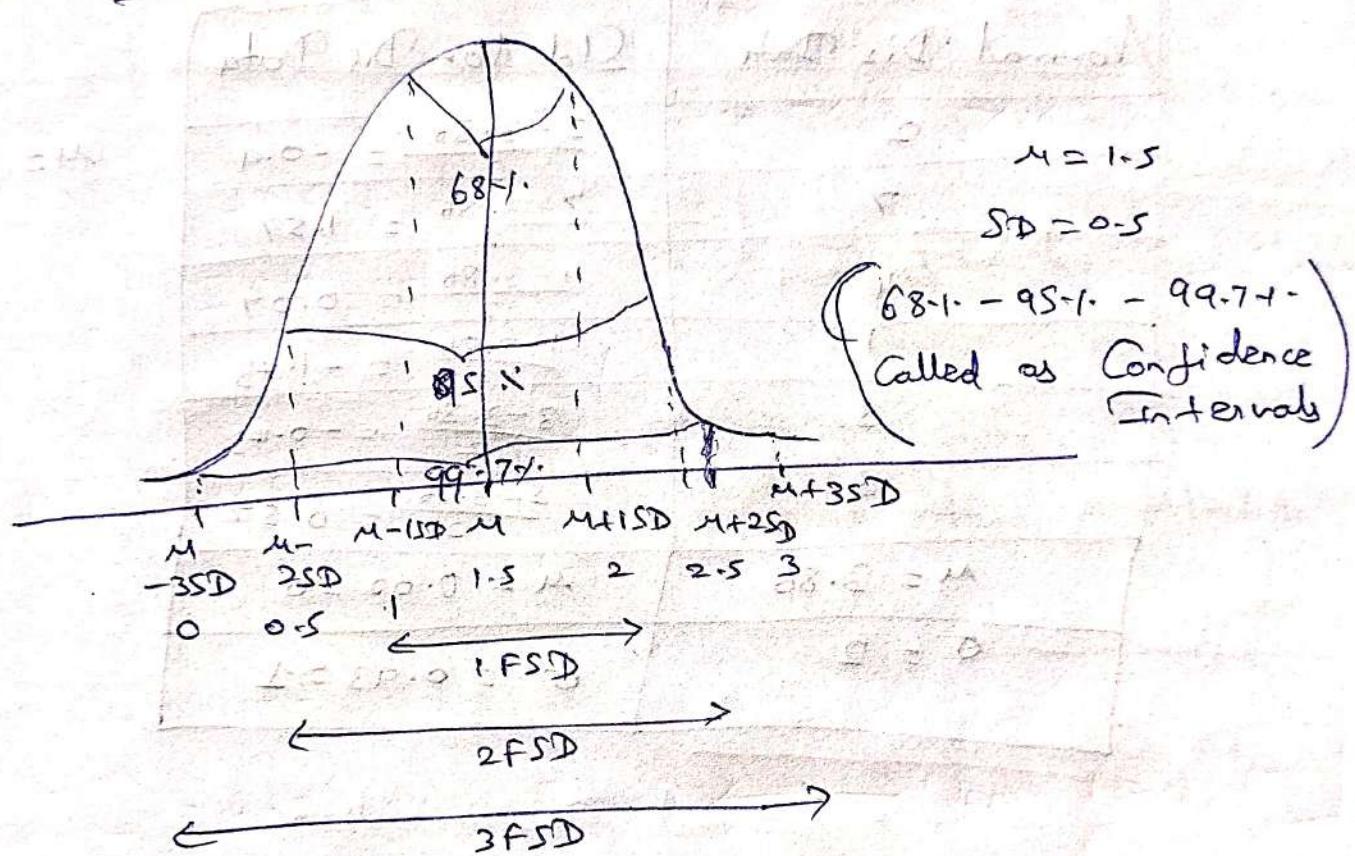
- 1) To understand data Patterns
- 2) To Summarize the data Easily
- 3) To Calculate Probabilities.
- 4) To Make Prediction & Decision
- 5) To choose right statistical test

There are 2 Category of Distribution

- 1) Continuous Distribution (Numerical)
- 2) Discrete Distribution (Categorical)

- 1) Normal Distribution
 - 2) Standard Normal Distribution
 - 3) Bernoulli Distribution
 - 4) Binomial Distribution
 - 5) Poisson Distribution
- Continuous Dist
(Numerical)
- Discrete Dist
(Categorical)

Normal Distribution:-



Mean = Median = Mode

Empirical Rules:-

- 68.3% of data will Present in 1st SD
- 95.4% of data will Present in 2nd SD
- 99.7% of data will Present in 3rd SD

27 Standard Normal Distribution

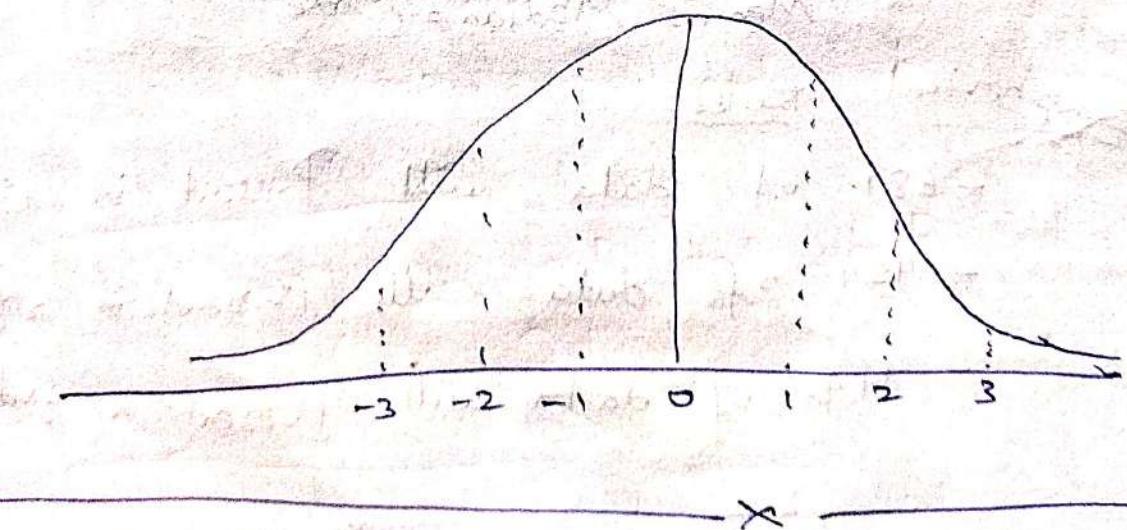
Always mean $\mu = 0$

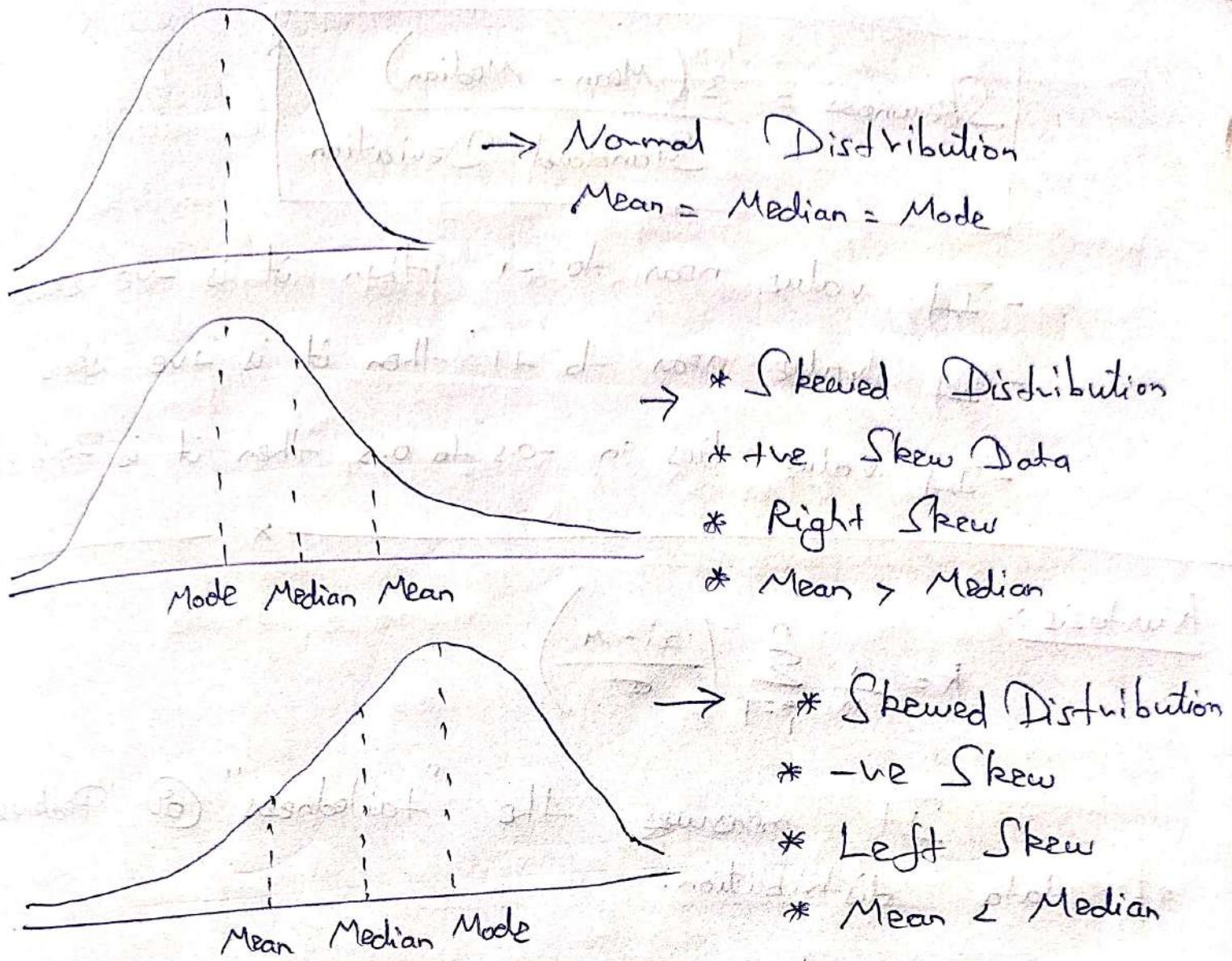
Always $SD = 1$

To convert Normal \rightarrow Standard Normal

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

Normal Dis Data	Std Non Dis Data	
2	$\frac{2-3.86}{2} = -0.9$	$\mu = 3.86$
7	$\frac{7-3.86}{2} = 1.57$	
4	$\frac{4-3.86}{2} = 0.07$	
1	$\frac{1-3.86}{2} = -1.43$	
3	$\frac{3-3.86}{2} = -0.43$	
5	$\frac{5-3.86}{2} = 0.57$	
$\mu = 3.86$	$\mu = 0.02 \approx 0$	
$\sigma = 2$	$\sigma = 0.93 = 1$	





4/11/2025

Positive Skew (Right skew):
Tail on the right side is longer.
Most data are on the left.

Negative Skew (Left skew):
- Tail on the left side is longer.
- Most data are on the right.

Zero Skew (Symmetric):
The data is evenly distributed around the mean (like a normal distribution).

$$\text{Skewness} = \frac{3(\text{Mean} - \text{Median})}{\text{Standard Deviation}}$$

- If value near -1 then it is -ve skew
- If value near $+1$ then it is +ve skew
- If value lies in -0.5 to 0.5 then it is Zero skew



Kurtosis :-

$$K = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{X}}{\sigma} \right)^4$$

If measures the "tailedness" or "Peakness" of data distribution.

Types of kurtosis

1) Mesokurtic ($K=3$)

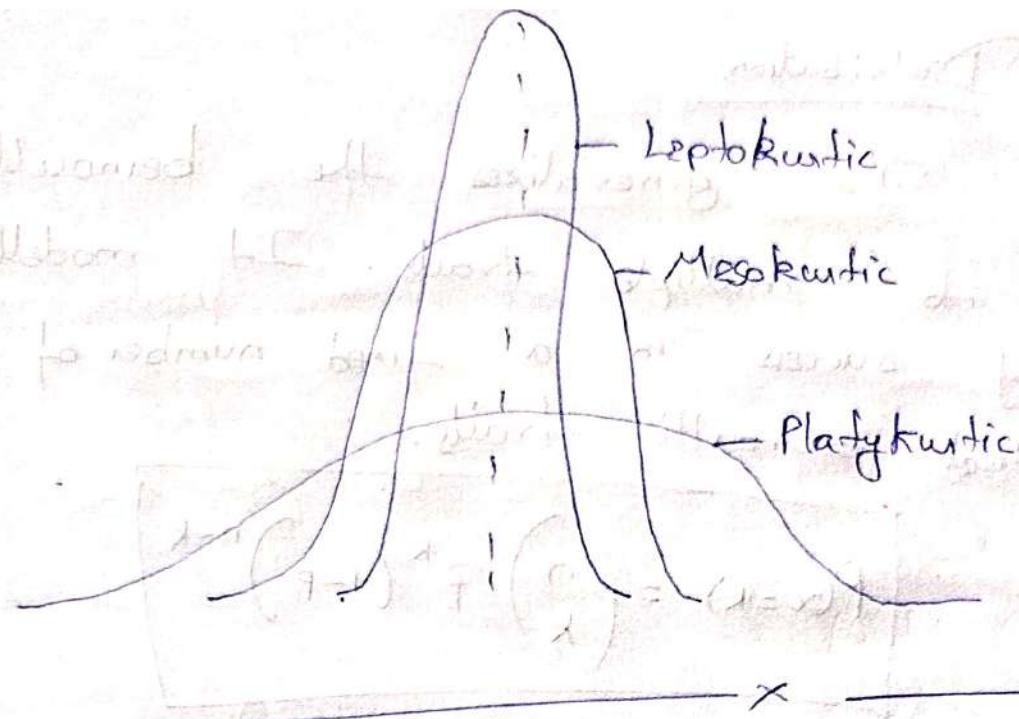
- * Normal Distribution
- * No outliers
- * Moderate tail and Peak

2) Lepokurtic ($K > 3$)

- * Heavy tails and sharp Peak
- * More Outliers

3) Platykurtic ($K < 3$) :-

- * Light tails and flat Peak
- * Fewer outliers



Bernoulli Distribution:

It is the simplest form of a discrete probability distribution and models a random experiment with exactly 2 outcomes.

Success denoted by = P

Failure denoted by = $1-P$

Total Probability = 1

The range is (0-1)

$$P(T) = \frac{1}{2} = 0.5$$

$$P(F) = \frac{1}{2} = 0.5$$



Binomial Distribution

It generalises the Bernoulli distribution to multiple trials. It models the number of success in a fixed number of independent and identical Bernoulli trials.

$$P(x=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

k = no. of success

n = no. of trials

p = probability of success

Poisson Distribution

It is used to model the number of events that occur in a fixed time interval.

(or) Space & occur independently the parameter λ represent the avg no of Event in the interval.

$$P(x=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

λ - Average no of Events

k - No of arguments

Interventional Statistics

Probability:

- It measures likelihood of an event.

Eg: Dice = $\{1, 2, 3, 4, 5, 6\}$

$$P(A) = \frac{\text{No of favourable outcomes}}{\text{Total no. of outcomes}}$$

$$P(B) = \frac{1}{6}$$

$$P(2, 4, 5) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} = 0.5$$

Toss the 2 coins: $\{HH, HT, TH, TT\}$

What is the P of getting only 1 H = $\frac{2}{4} = \frac{1}{2}$

$$\text{Both tail} = \frac{1}{4}$$

There are two rules in Probability:-

⇒ Addition Rule (OR)

⇒ Multiplication Rule (AND)

Addition Rule:-

1) Mutually Exclusive Event

2) Non-Mutual Exclusive Event

5/11/2025

Unit 2 Probability

⇒ Mutual Exclusive Events:-

The different events won't occur at the same time.

Ex: If you toss the coin, what is the probability of landing on head or tail.

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

$$\begin{aligned} P(H \text{ or } T) &= P(H) + P(T) \\ &= \frac{1}{2} + \frac{1}{2} \end{aligned}$$

$$P(H \text{ or } T) = 1$$

⇒ Non-Mutual Exclusive Events:-

Multiple Events can occur at the same time.

Ex: Picking the cards from deck cards.

What is the Probability of getting Jack or King?

$$P(J \text{ or } K) = P(J) + P(K) - P(J \cap K)$$

$$= \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

$$= \frac{16}{52}$$

$$P(J \text{ or } K) = 0.31$$

∴ 31 - 1.

Multiplication Rule :-

1) Independent Events-

Here all the values have the same priority after n number trials also.

(or) 1 Event don't depend on another event.

Eg: 1st toss of Coin

$$P(H) = \frac{1}{2} \quad n^{\text{th}} \text{ Toss of coin}$$

2nd Toss of Coin

$$P(H) = \frac{1}{2}$$

$$P(H) = \frac{1}{2}$$

Eg: What is the probability of Dice rolling and getting a 5's and 4's?

$$P(A \text{ and } B) = P(A) \times P(B)$$

$$P(5 \text{ and } 4) = P(5) \times P(4)$$

$$= \frac{1}{6} \times \frac{1}{6}$$

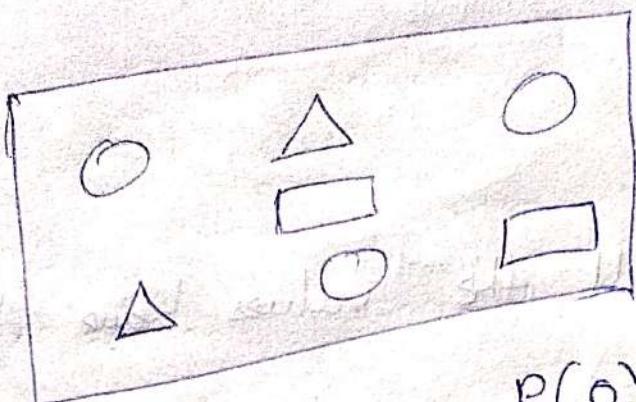
$$= \frac{1}{36}$$

$$P(5 \text{ and } 4) = 0.027 \rightarrow 2.7\%$$

2) Dependent Events-

Previous

Present Event is depend on the Event.



\Rightarrow Total number of shapes = 7

$$1^{\text{st}} \text{ time} \rightarrow P(O) = \frac{3}{7} = 0.428$$

$$2^{\text{nd}} \text{ time} \rightarrow P(\Delta) = \frac{2}{6} = 0.333$$

$$3^{\text{rd}} \text{ time} \rightarrow P(\square) = \frac{2}{5} = 0.4$$

$$4^{\text{th}} \text{ time} \rightarrow P(O) = \frac{2}{4} = 0.5$$

Ex: From a deck of cards what is the probability of getting a king and then 8?

$$P(A \text{ and } B) = P(A) \times P(B|A)$$

$$P(K \text{ and } 8) = P(K) \times P(8|K)$$

$$= \frac{4}{52} \times \frac{4}{51}$$

$$P(K \text{ and } 8) = 0.07 \rightarrow 7\% \text{ chance}$$

Permutation

Permutation

Combination

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

{ Dosa, Idly, Vada, Puri }

DI	ID	VD	PD
DV	IV	VI	PI
DP	IP	VP	PV

$$\therefore {}^n P_r = \frac{4!}{(4-2)!} = 12,$$

$n \rightarrow$ no of items

$r \rightarrow$ no of ways

Permutation refers to the different ways in which a set of items can be arranged in order and in permutation the order of the items matters but not items.

2) Combinations:-

{ Dosa, Idly, Vada, Puri }

DI	6
DV	
DP	
IV	
IP	
VP	

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

$$= \frac{4!}{(4-2)! 2!}$$

$$= \frac{4 \times 3 \times 2 \times 1}{(2 \times 1) \cdot (2 \times 1)}$$

$$= 6$$

It refers to the different way of selecting item from a set where the order of selection doesn't matter but items should not repeat.

6/11/2025

Hypothesis Testing

$H_0 / H_N \rightarrow$ Null Hypothesis (True Statement)

$H_1 / H_A \rightarrow$ Alternative Hypothesis

H_0 - The avg height of Indian is 5.7

H_{A1} - No, the avg height of Indian is not 5.7

H_0 - True Statement

H_{A1} - It will helps to find either we have to accept or reject H_0 .

P-value

If P value $\leq \alpha$ we can reject H_0

If $P > \alpha$ we can accept H_0

α - Significance value

$$\delta = 1 - \alpha$$

CI (Confidence Interval)

68%

95%

99.7%

If CI is 68% $\rightarrow \delta = 0.32$

If CI is 95% $\rightarrow \delta = 0.05$

If CI is 99.7% $\rightarrow \delta = 0.003$

P. Value formula:-

$$Z = \frac{\hat{P} - P_0}{\sqrt{\frac{P_0(1-P_0)}{n}}}$$

\hat{P} = Sample Proportion

P_0 = Assumed Population Proportion in the H_0

n = Sample Size

$P = 0.35, \delta = 0.05 \rightarrow$ accept

$P = 0.03, \delta = 0.05 \rightarrow$ reject

- Hypothesis Testing is a framework

for making inferences about data and models
in Machine Learning.

- It helps in model evaluation, feature selection, assumption validation and ensuring the fast robustness & reliability of conclusions drawn from models.

Type I and Type II Error:-

R H_0 is True, D H_0 is True ✓

R H_0 is True, D H_0 is False (Type I Error)

R H_0 is False, D H_0 is True (Type II Error)

R H_0 is False, D H_0 is False ✓

R - Reality D - Decision

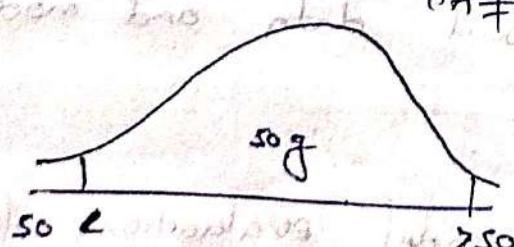
Type I - If you failed to accept H_0

Type II - If you failed to reject H_0 (Most dangerous)

One tail test and two tail test

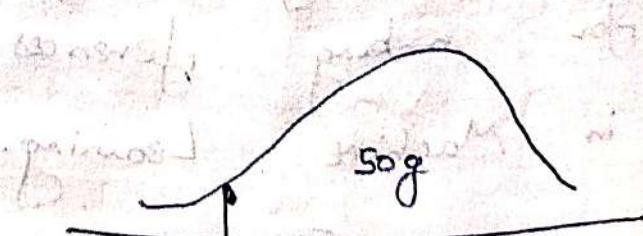
$H_0 \rightarrow$ The chips packet weight is 50g

$H_A \rightarrow$ The chips packet weight is not 50g.



H_0 (Two tail)

$$\begin{aligned} H_0 &= 50g \\ H_A &\neq 50g \end{aligned}$$



H_A (One tail)

Z test & T test:

The avg age of an college stu is 24 years with $SD = 1.5$, Sample of 36 Students, The mean is 25 years with 95% confidence Interval. Do the age will vary or not?

Sol:

Given: Population mean $\mu = 24$

Population SD $\sigma = 1.5$

Sample mean $\bar{x} = 25$

Confidence Interval = 95%.

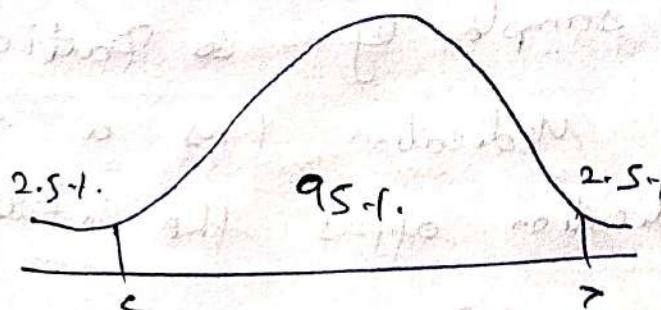
$$\alpha = 1 - 95\% = 0.05$$

If they given Population SD - "Z test"

If they given Sample SD - "T test"

Null Hypothesis $H_0 \rightarrow$ Avg age is 24

Alternate Hypothesis $H_A \rightarrow$ Avg age is not 24



$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$= \frac{25-24}{\frac{1.5}{\sqrt{36}}} = \frac{1}{(1.5/6)} = 4$$

$$\boxed{Z=4}$$

In Z table value for 4 with α of 0.05

$$\therefore 0.99997$$

$$\boxed{\text{AUC} = 1}$$

$$\therefore 1 - 0.99997 = 0.00003$$

$$= \frac{0.00003}{2}$$

$$\boxed{P = 0.000015}$$

$$\boxed{\delta = 0.05}$$

$$\therefore P < \delta$$

\therefore So we can Reject H_0

In the population the avg IQ is $100^{\prime \prime}$

$SD = 15$. Researcher want to test a new medication see if there is any effect on intelligence or no effect at all a sample of 36 participants who have taken the medication has a mean of 140 did the medication offer the intelligence of 95%?

$$M = 100 \quad n = 30 \quad \bar{x} = 140$$

$$CI = 95\% \quad d = 0.05$$

$$Z = \frac{\bar{x} - M}{\sigma / \sqrt{n}} = \frac{140 - 100}{15 / \sqrt{30}} = 14$$

Z table check value for 14 with $d=0.05$

$$Z = 14 \rightarrow 1$$

$$P = 1 - 1$$

$$P = \frac{0}{2} = 0$$

$P < d = 0.05$ we reject

H_0 = The avg IQ is 100

H_A = The avg IQ is not 100

$$M = 130 \quad n = 30 \quad \bar{x} = 140$$

$$CI = 95\% \quad d = 0.05$$

$$Z = \frac{\bar{x} - M}{\sigma / \sqrt{n}} = \frac{140 - 130}{15 / \sqrt{30}}$$

$$Z = \frac{10}{2.7}$$

$$Z = 3.66 \quad d = 0.05$$

X

7/11/2025

Ex:

$$\mu = 100$$

$$n = 30$$

$$\bar{x} = 140$$

$$S = 20$$

$$\alpha = 0.05$$

Sample SD, S_0 & t_{test}

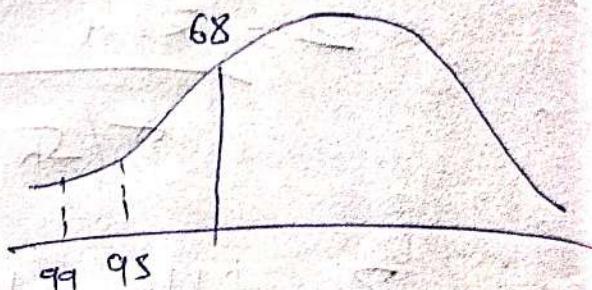
$$t = \frac{\bar{x} - \mu}{S/\sqrt{n}}$$

$$= \frac{140 - 100}{\frac{20}{\sqrt{30}}}$$

$$= \frac{40}{20/\sqrt{30}}$$

$$= \frac{40}{3.6514}$$

$$t = 10.95$$



$$\text{DOF} = n - 1$$

$$230 -$$

$$229$$

$t > t_{\text{table}}$ value with DOF = 29 & $\alpha = 0.05$

$H_0 \rightarrow \text{Avg IQ} = 100$ }
IQ $\neq 100$ } Two tail

\therefore Table value = 2.045

$\therefore t \text{ value} > \text{Table value}$

$$10.95 > 2.045$$

\therefore we can reject,

Ex:

Credit Card Launch:

$$n = 140 \quad \bar{x} = 1990 \quad \sigma = 2500 \quad s = 2833$$

$$CI = 95\% \quad Z_{\alpha/2} = 1.96$$

We have to calculate the range interval.

Default always the CI will be 95%.

$$CI = \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Lower bound

$$= \bar{x} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$= 1990 - \frac{20.05}{2} \frac{2500}{\sqrt{140}}$$

$$= 1990 - 20.025 (211.29)$$

$$= 1990 - (1.96) (211.29)$$

$$= 1990 - (414.12)$$

$$= 1575$$

Upper bound

$$\bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$= 1990 + \frac{20.05}{2} \left(\frac{2500}{\sqrt{140}} \right)$$

$$= 1990 + (1.96) (211.29)$$

$$= 1990 + 414.12$$

$$= 2404$$

i. The avg balance they are going to maintain after fully fledged launch is $[1575 \quad 2404]$

Ex: On a quant test of a CAT Exam of a sample of 25 test takers as a sample mean of 520 with sample standard deviation of 80, construct 95.1% confidence interval about the mean.

Sol:

$$n = 25 \quad \bar{x} = 520 \quad S = 80 \quad CI = 95\% \quad d.f = 24$$

$$CI = \bar{x} \pm t_{0.025} \frac{S}{\sqrt{n}} \quad DOF = n - 1 = 24$$

Lower Band

$$CI = 520 - t_{0.025} \frac{80}{\sqrt{25}}$$

$$= 520 - t_{0.025}(16)$$

$$= 520 - (2.064)(16)$$

$$CI = 486.976$$

\therefore Range [487

Upper Band

$$CI = 520 + t_{0.025} \frac{80}{\sqrt{25}}$$

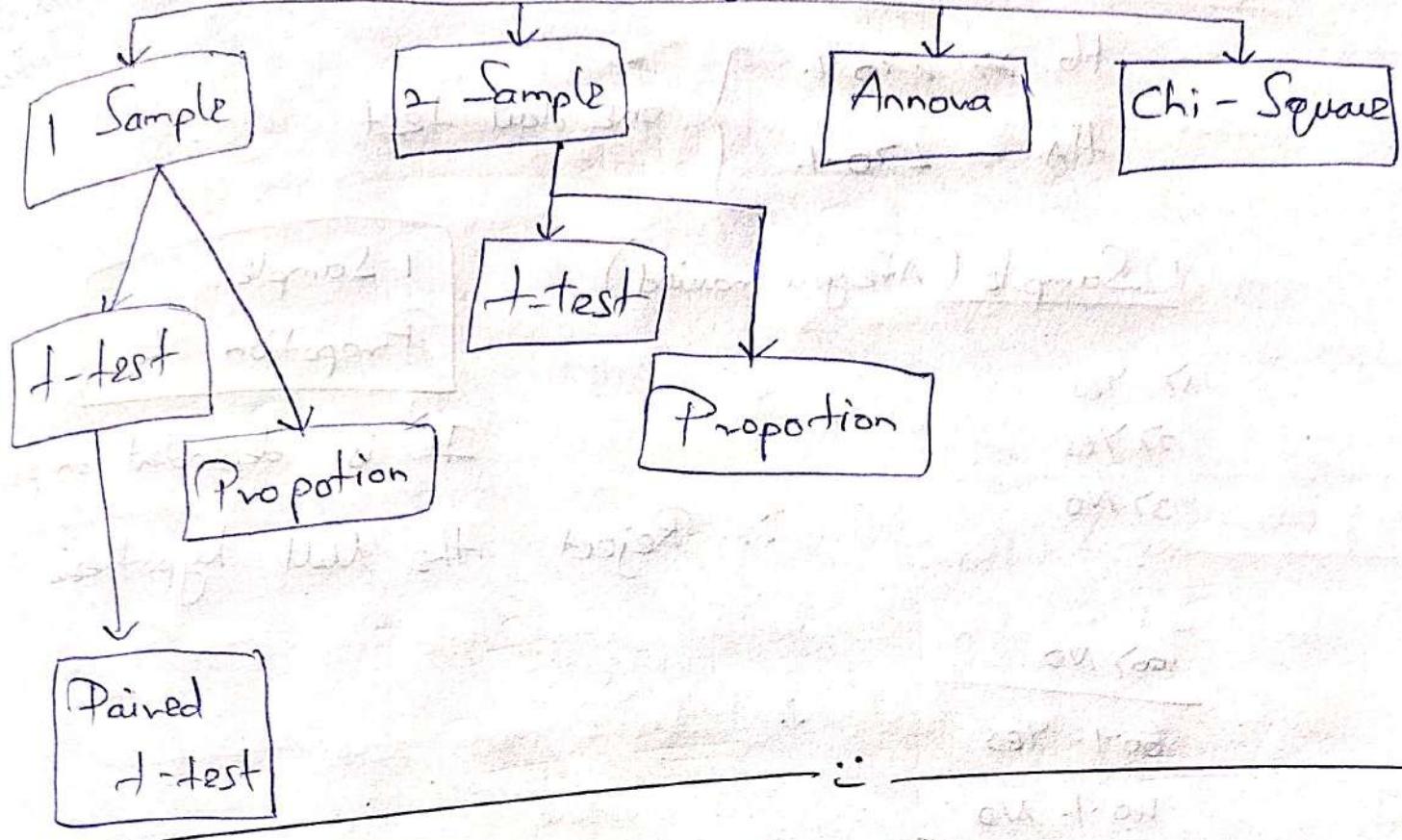
$$= 520 + t_{0.025}(16)$$

$$= 520 + 2.064(16)$$

$$CI = 553$$

553]

Statistical tests



10/11/2025

H_0 - The avg sal of IT Emp in BLR is 27 k
 H_A - The avg sal of IT Emp in BLR is not 27 k.

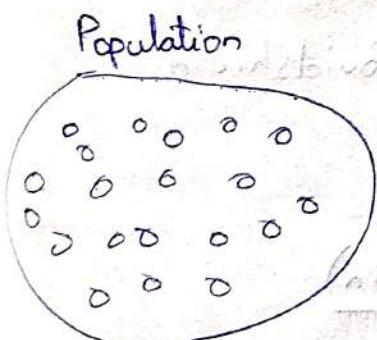
1 Sample t-test → depends on population

$H_0 \rightarrow \text{sal} = 27\text{ k}$

$H_A \rightarrow \text{sal} \neq 27\text{ k}$

Sal < 27 k \rightarrow

Two tail test



Sample 1

1) 28 k

2) 30 k

3) 18 k

⋮

⋮

($\text{avg} > 25\text{ k}$)

Avg = 32 k

Reject Null Hypothesis

\therefore The avg sal of IT Emp IN BLR is 32 k

H_0 - More than 70% of People are married in India
 H_A - No more than 70% of People are ~~not~~ married in India.

$$\begin{array}{l} H_0 \Rightarrow > 70\% \\ H_A \Rightarrow \leq 70\% \end{array} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{one tail test}$$

1 Sample (Are you married)

1) Yes

2) Yes

3) No

:

100) No

60% Yes

40% No

\therefore Reject the null hypothesis

1 Sample
Proportion Test

\downarrow is dependent on population

2 Sample t-test :- (Independent on Population) :-

H_0 : Covidshield is better than Co-vaxin

H_A : No, Co-vaxin is better than Covidshield

How many hrs it take to react

Sample 1 (covidshield)

- 1) 2hr
- 2) 3hr

:

50) 1 hr

3.5hr

Sample 2 (co-vaxin)

- 1) 4hr
- 2) 5hr

:

50) 2hr

7.2hr

∴ We have to accept H_0 , because we got 3.5 hrs reaction time for covidshield & 7.2 hrs for co-vaxin

2-Sample Proportion : Independent on Population:-

H_0 - New Beauty treatment is better than old one

H_A - Old Beauty treatment is better than new one

Are you satisfied with the treatment

New (S_1)

1) Yes

2) No

:

50) Yes

50% Yes

50% No

Old (S_2)

1) No

2) Yes

:

50) No

50% Yes

50% No

∴ We can reject H_0

Bcz 80% of people like old

1 Sample Paired t-test

$H_0 \rightarrow$ By joining new weight loss program you can see significant difference in your weight.

$H_A \rightarrow$ By joining new weight loss program, you can't see significant difference in your weight.

Sample I

<u>Before</u>	<u>After</u>
1) 70	1) 65 ✓
2) 80	2) 70 ✓
:	:
10) 100	10) 102 ✗

we can accept the H_0 ,
bcz majority got weight
loss.

11/11/2025

H_0 - Your batch students can't able to score $> 90\text{ M}$

H_A - No my batch students can able to score $> 90\text{ M}$

B ₁	B ₂	B ₃	B ₄	B ₅
63	75	92	47	81
62	88	98	58	75
58	83	88	63	78
73	65	69	70	85
80	73	95	68	70
65	78	92	59	80

ANOVA

\therefore If one sample proved means we can reject the H_0 .

\therefore We reject the null hypothesis

In the 2000 Indian Census the age of the individual in a small town were found to be following.

In the year 2000

Less than 18	18 - 35	> 35
204	304	56

In 2010 Age of $n=500$ individuals were sampled below are the results.

In year 2010

Less than 18	18 - 35	> 35
121	288	91

Using $\alpha = 0.05$ would you calculate the population distribution of ages has changed in the last 10 years?

Year	< 18	18 - 35	> 35	
2000	204	304	56	
2010	100	150	250	→ observed
	121	288	91	→ observed

H_0 — 2010 census ratio same as 2000

H_A — No, 2010 census is not same as 2000

Chi-Square Test:-

$\chi^2 >$ chi-square table value with Degree of freedom means you can reject Else accept.

$$\text{Degree of freedom} = 3-1 = 2 - \text{Dof}$$

$$\alpha = 0.05$$

$\therefore \chi^2 > 5.991$ then reject H_0

$$\chi^2 = \sum_{i=1}^r \frac{(f_o - f_e)^2}{f_e}$$

f_o - observed value

f_e - expected value

$$\chi^2 = \frac{(121-100)^2}{100} + \frac{(288-250)^2}{250} + \frac{(91-250)^2}{250}$$

$$\chi^2 = 232.494 > 5.991$$

\therefore So we can reject Null Hypothesis H_0

Correlation:-

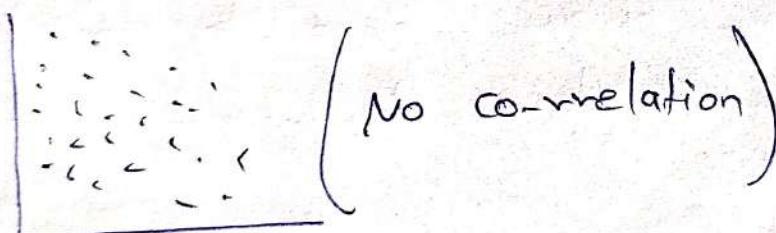
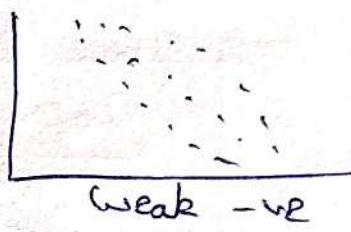
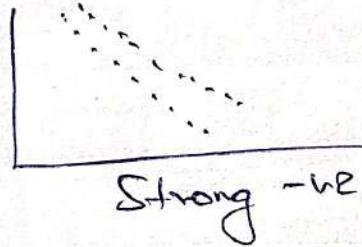
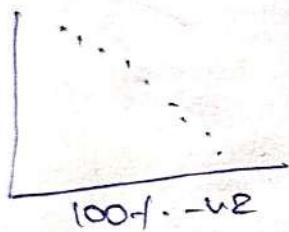
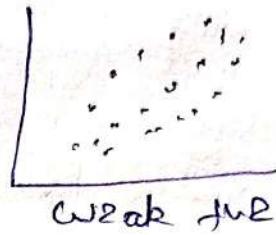
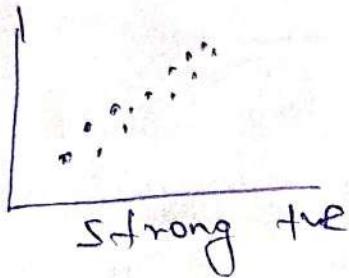
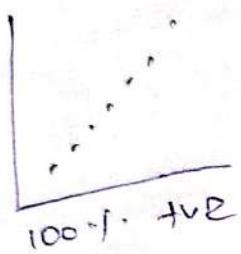
To check the relationship b/w two numerical features/column we use correlation.

There are 3 types of Correlation

1) +ve Corre \rightarrow $x \uparrow y \uparrow$ \rightarrow Exp T, Sal \uparrow

2) -ve Corre \rightarrow $x \uparrow y \downarrow$ \rightarrow Weight \uparrow , MPG \downarrow

3) No Corre \rightarrow $x \uparrow y \uparrow$ \rightarrow weight, Sal



(No correlation)

There are two types of correlation formula

1) Pearson - Correlation

2) Spearman - Correlation

- The correlation value ranges from -1 to +1
- If correlation value near do +1 is a +ve-Corr
- If correlation value near do -1 is -ve-Corr
- If value near do Zero - No-Corr.