

31/10/25

Measure of dispersion:

Variance:

Ex:

	person 1
Mon	7:00 AM
Tue	7:15 AM
Wed	7:30 AM
Thur	7:45 AM
Fri	8:00 AM
Sat	9:00 AM - 8:00 AM.
	7:15 - 15 mm
	7:30 - 30 m
	7:45 - 15 m

person 2

8:00 AM
9:00 AM
11:00 AM
7:00 AM
10:00 AM
? 9-10.

7 AM - 2 hr
11 AM - 1 hr

Variance high prediction are low.

1. Variance
2. Standard Deviation.
3. Range

Variance:

1. population variance (σ^2)

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

2. Sample variance (S^2)

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

"n-1" we call it as degree of freedom.

Ex: Calculate variance on { 1, 2, 2, 3, 4, 5 }.

$$\begin{aligned} \text{population } \sigma^2 &= \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} \\ &= \frac{1+2+2+3+4+5}{6} \\ &= \frac{17}{6} \\ &= \frac{(1-2.8)^2 + (2-2.8)^2 + (3-2.8)^2 + (4-2.8)^2 + (5-2.8)^2}{6} \\ &= \frac{3.36 + 0.69 + 0.69 + 0.02 + 1.36 + 4.69}{6} \\ &= \frac{10.81}{6} = 6.90 \end{aligned}$$

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

$$\sigma^2 = \sqrt{1.8 \text{ km}^2} \text{ (Variance)}$$

$$\sigma = 1.34 \text{ km (Standard deviation)}$$

$$\text{Standard Deviation} = \sqrt{\text{Variance}}$$

Sample SD

$$\text{population SD} \\ \sigma = \sqrt{\sigma^2}$$

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

$$S = \sqrt{2.16}$$

$$S = 1.46$$

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

$$\sigma = 1.34$$

$$\text{Range} = \text{Maximum} - \text{Minimum}$$

$$= 5 - 1$$

$$\text{Range} = 4$$

percentile and Quartile :

percentile is the value below which certain percentage of observation will ~~lie~~ lie.

Ex :

= { 1, 1, 2, 3, 4, 5, 5, 6, 7, 7, 8 }.

How much % of data will come below 6 ?

$$a) \text{ percentile rank of } x = \frac{\# \text{ no of value below } x}{N} \times 100.$$

$$= \frac{7}{11} \times 100$$

$$= 63.63$$

Observation of data value is < 6 .

Quartile :

Quartile helps to find the value which is present at the given percentile rank.

Ex :

=

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

Which value is present at 25% ?

$$\text{value} = \frac{\text{percentile}}{100} \times n+1$$

$$= \frac{25}{100} \times 12$$

$$= 3 \rightarrow \text{Index}$$

Value = 2

Calculating for 90%.

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

$= 10.8 \xrightarrow{10}$ This value called Index.

$$= 7$$

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.

-4.5 15.5
[Lower fence upper fence]

Anything < -4.5 & > 15.5 is a

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

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

IQR (Inter Quartile Range)

$$IQR = Q_3 - Q_1$$

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

$$Q_1 = \frac{25}{100} \times 18$$

$$= 4.5 = \text{Index}$$

$$Q_1 = 3$$

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

$$= 13.5 = \text{Index}$$

$$Q_3 = 8$$

$$IQR = Q_3 - Q_1$$

$$= 8 - 3$$

$$IQR = 5$$

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

$$= -4.5$$

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

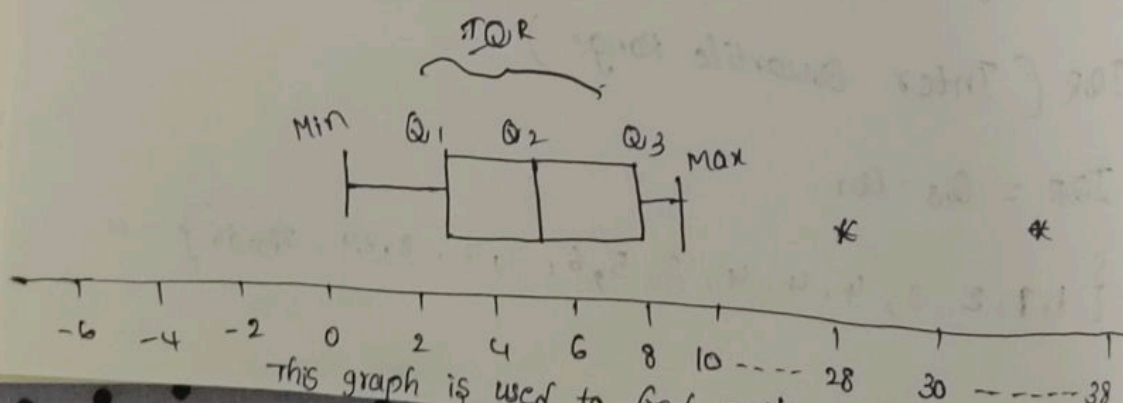
$$= 15.5$$

Minimum = 1, Median = 5, Maximum = 9

$$Q_1 = 3$$

$$Q_3 = 8$$

Box plot:



3/11/25

Different types of Distributions:

1. To understand data patterns.
2. To summarize the data easily.
3. To calculate the probabilities.
4. To make prediction and decision.
5. To choose right statistical test.

There are two category of distribution.

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

1. Normal Distribution.
2. Standard normal Distribution.

} Continuous Distribution

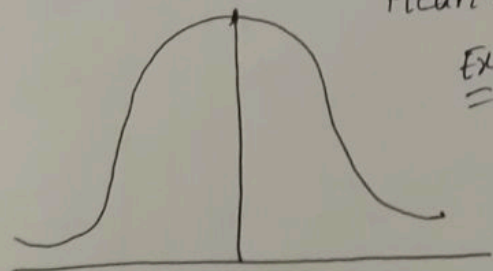
3. Bernoulli Distribution.
4. Binomial Distribution.
5. Poisson Distribution.

} Categorical Distribution.

Normal Distribution:

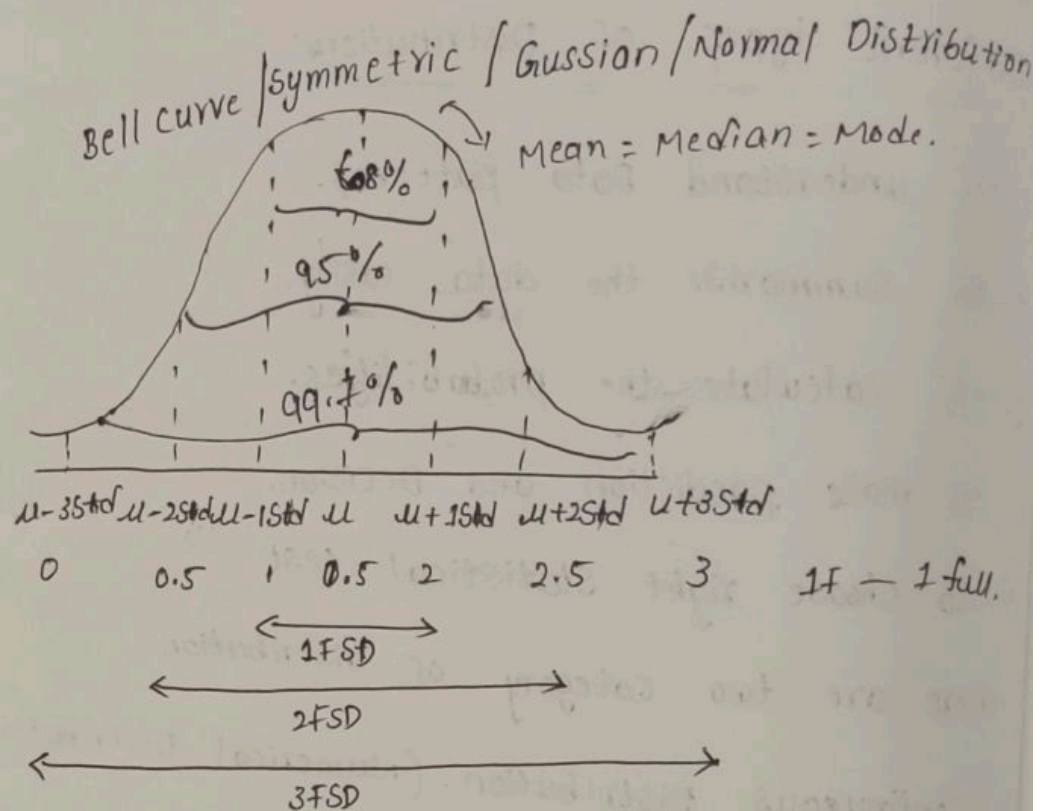
Mean = Median = Mode.

Ex!



1. Temp
2. Speed
3. Strength level

The Normal Distribution follows Rules of empirical.



→ 68% - 95% - 99.7% Called as Confidence Intervals.

Empirical Rules:

68% of data will present in 1SD.

95% of data will present in 2SD.

99.7% of data will present in 3SD.

Standard Normal Distribution (SND):

$$\mu = 0$$

$$\sigma = 1$$

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

Normal Distribution data

Standard Normal Distribution data.

2

-0.9

7

1.57

5

0.57

4

0.07

1

-1.43

3

-0.43

5

0.57

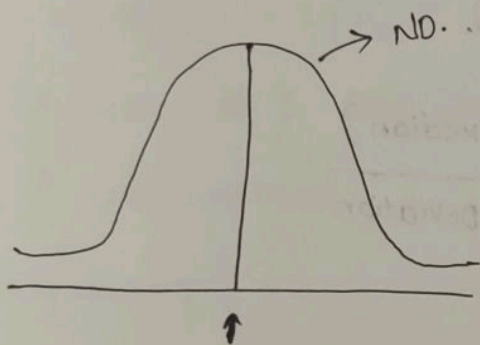
$$\mu = 3.86$$

$$\mu = 0.02$$

$$\sigma = 2.$$

$$\sigma = 0.93 = 1$$

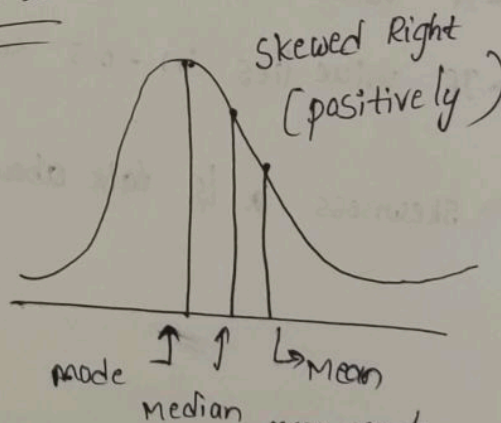
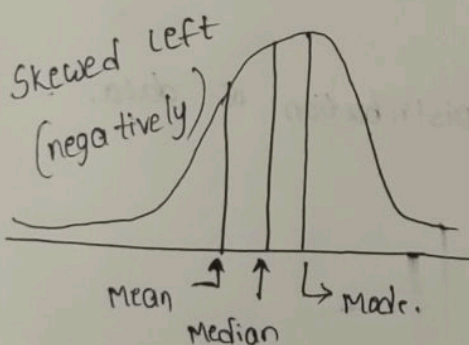
Normal Distribution:



Mean = Median = Mode.

Symmetric

Skewed Data



4/11/25

Different types of Skews:

1. positive skew (Right side):

Tail on the right side is longer most data are on the left.

2. Negative skew (left side):

Tail on the left side is longer most data are on the right.

3. Zero skew (Symmetric):

The data is evenly distributed around the mean (like a normal distribution).

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

It is only for
Individual columns

→ If value near to -1 then it is -ve skew.

→ If value near to +1 then it is +ve skew.

→ If value lies in -0.5 to 0.5 then it is zero skew.

→ skewness is talk about distribution of data.

Kurtosis :

$$k = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma} \right)^4$$

It measures the "tailedness" or "peakness" of data distribution.

Types of kurtosis :

1. Meso kurtic ($k=3$).

- * Normal Distribution.
- * No outliers.
- * Moderate tail and peak.

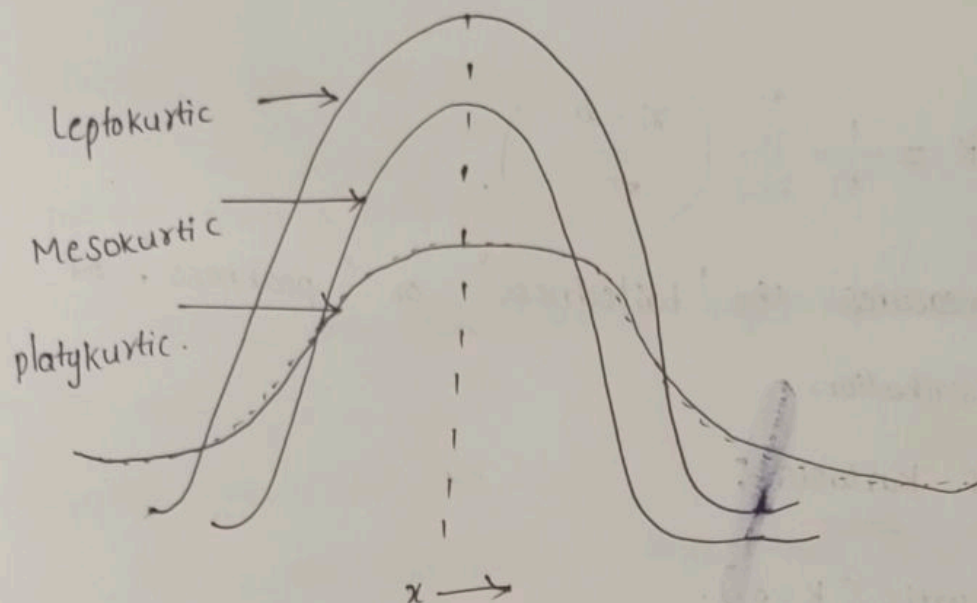
2. Lepo kurtic ($k>3$).

- * Heavy tails and sharp peak.
- * More outliers.

3. platy kurtic ($k<3$).

- * Light tails and flat peak.
- * Fewer outliers.





Discrete Distribution:

Bernoulli Distribution:

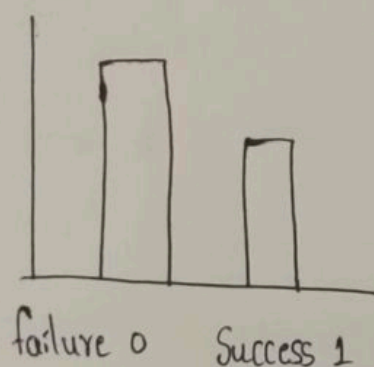
It is simplest form of discrete probability distribution. and models a random experiment with exactly 2 outcomes.

Success denoted by $= p$

failure denoted by $= 1-p$.

→ Total probability is 1.

→ The range of probability is 0-1.



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

$$p(H) = \frac{1}{2} = 0.5$$

Binomial Distribution:

It generalizes 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 and occur independently. The parameter λ represents the avg number of event in the interval.

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

$k!$ = factorial of k .

λ = Avg no. of events.

k = No. of (Arguments.) occurrences.

e = Euler's number.

Inferential Statistics :

1. probability :

It measures likelihood of an event.

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

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

$$p(3) = \frac{1}{6}$$

$$p(2, 4, 5) = p(2) + p(4) + p(5)$$

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

$$= 0.5$$

Toss the 2 coins.

$\{HH, HT, TH, TT\}$.

1) What is the probability of getting only 1 H?

$$A) = \frac{2}{4_2} = \frac{1}{2}.$$

2) What is the probability of getting both Tails?

$$\frac{1}{4}$$

There are 2 rules in probability :

- 1) Addition Rule (OR)
- 2) Multiplication Rule (AND)

Addition Rule :

1. Mutual Exclusive Events.
2. Non Mutual Exclusive Event.

5/11/25

Mutual Exclusive Event :

The different events won't ~~never~~ ^{occur} at the same time is called Mutual Exclusive Event.

Ex: If you toss the coin what is the probability of landing on heads or tails.

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

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

$$P(H \text{ or } T) = \frac{1}{2} + \frac{1}{2} = 1.$$

Non Mutual Exclusive Event :

Here multiple events can occur at the same time. is called Non Mutual Exclusive Event.

Ex: picking the cards from deck cards

What is the probability of getting Jack or heart?

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

$$P(J \text{ or } H) = \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

$$P(J \text{ or } H) = 0.07 + 0.25 - 0.01$$

$$P(J \text{ or } H) = 0.33 = 0.31$$

Multiplication Rule:

1. Independent Event:

Here all the values have the same probability after n. number trials also (or) 1 Event don't depend on another event.

1st toss the coin

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

7th toss the coin

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

Ex: What is the probability of Dice rolling and getting a 5 and then 4?

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

$$p(5 \text{ and } 4) = p(5) * p(4)$$

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

$$= \frac{1}{36} = 0.027$$

2.7 % is getting.

2) Dependent Event:

present event depend on the previous event.

1st time

$$p(O) = \frac{3}{7} = 0.43$$

2nd time

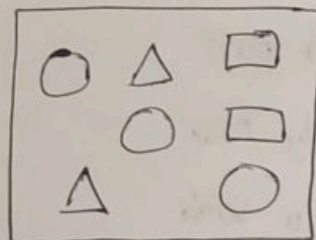
$$p(\Delta) = \frac{2}{6} = 0.3$$

3rd time

$$p(\square) = \frac{2}{5} = 0.4$$

4th time

$$p(O) = \frac{2}{4} = 0.5$$



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

$$p(A \text{ and } B) = p(A) * p(B/A)$$

$$p(K \text{ and } 8) = \frac{4}{52} * \frac{4}{51} = 0.07 * 0.07$$

$$p(K \text{ and } 8) = 0.07 = 7\%$$

permutation and Combination :

permutation means the value should be repeated but order shouldn't be repeated.

Ex: Dosa, Idly, Vada, puri

1. Dosa Idly

7. Vada Dosa

2. Dosa vada

8. Vada Idly

3. Dosa puri

9. vada puri

4. Idly Dosa

10. puri Dosa

5. Idly vada

11. puri Idly

6. Idly puri

12. puri Vada.

$${}^n P_r = \frac{n!}{(n-r)!} = \frac{4!}{(4-2)!} = \frac{4!}{2!} = \frac{4 \times 3 \times \cancel{2} \times 1}{\cancel{2} \times 1} = 12$$

n = no. of items

r = no. of ways.

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

Combination :

{ Dosa, Idly, vada, puri }

1. Dosa Idly

2. Dosa Vada

3. Dosa puri

4. Idly Vada

5. Idly puri

6. Vada puri

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

$$= \frac{4!}{2!2!}$$

$$= \frac{4 \times 3 \times 2!}{2!2!}$$

→ 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.

6/11/25

Hypothesis testing:

$H_0/H_N \rightarrow$ Null hypothesis (True Statements)

$H_1/H_A \rightarrow$ Alternative hypothesis.

Ex:

H_0 : The Avg height of Indian man is 5.7.

H_A : No, The Avg height of Indian man is not 5.7.

H_0 : True Statement.

H_A : It will help to find either we have to accept (or) reject H_0 .

p-value :

If $p < \alpha$ we can reject H_0 .

If $p > \alpha$ we have to accept the H_0 .

→ " α " is called as significance value.

α - Significance value :

$$\alpha = 1 - CI$$

→ If CI is 68% what is the α is "0.32".

→ If CI is 95% then α is "0.05".

→ If CI is 99.7% then α is "0.03".

CI = Confidence Interval.

P-value formula :

$$Z = \frac{\hat{p} - p_0}{\frac{p_0(1-p_0)}{n}}$$

\hat{p} = Sample proportion.

p_0 = Assumed population proportion in the H_0 .

n = Sample size.

$$p = 0.03$$

$$\alpha = 0.05$$

$p < \alpha$ reject H_0 .

→ Hypothesis testing is a framework for making inference about data and models in Machine Learning.

→ It helps in model evaluation feature selection assumption validation and ensuring the robustness and

reliability of conclusions drawn from the models.

Type I and Type II Error:

$R \rightarrow$ Reality.

$D \rightarrow$ Decision.

$R H_0$ is True $D H_0$ is True.

$R H_0$ is True $D H_0$ is false \rightarrow Type I Error.

$R H_0$ is false $D H_0$ is True \rightarrow Type II Error.

$R H_0$ is false $D H_0$ is false.

If you failed to accept H_0 is type I error.

If you failed to reject H_0 is type II Error.

and also most dangerous.

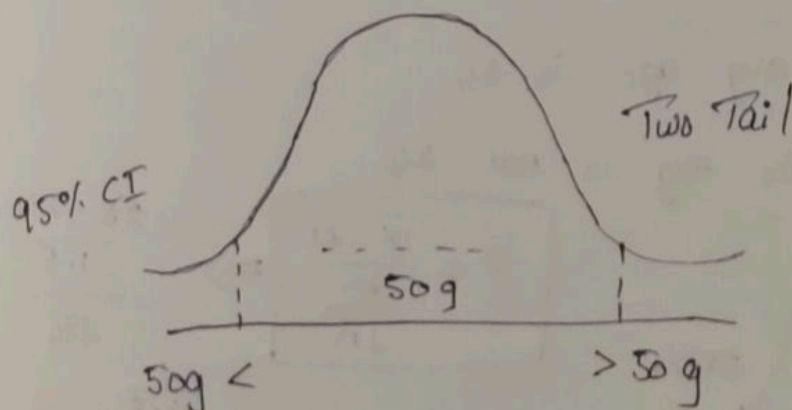
One Tail test and Two Tail test:

H_0 :- The chips packet weight is 50g.

H_A :- The chips packet weight is not 50g.

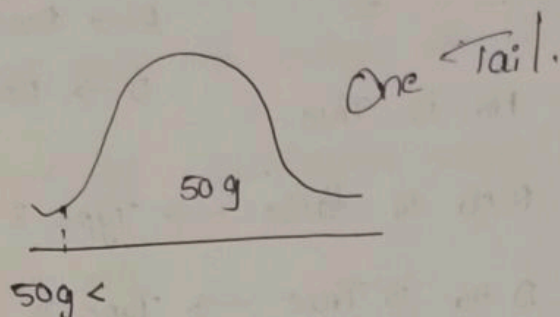
$$H_0 = 50g$$

$$H_A \neq 50g$$



→ H_0 :- The chips packet weight is more than 50g.

H_A :- No The chips packet weight is not more than 50g.



Z Test and T Test :

→ The Average age of college student is 24 years with the standard deviation 1.5. Sample of 36 students the mean is 25 years with 95% confidence interval. Do the age will vary or not.

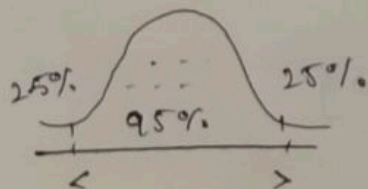
Z-Test : If they given population standard deviation go with z test.

T-test : If they give Sample standard deviation go with T test.

$$\mu = 25, \sigma = 1.5, n = 36, \bar{x} = 25, CI = 95\%, \alpha = 0.05$$

H_0 : The avg age is 24

H_A : No The Avg is not 24



$$z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \Rightarrow \frac{25 - 24}{\frac{1.5}{\sqrt{36}}}$$

$$z \Rightarrow \frac{1}{\frac{1.5}{6}} = 4$$

z table check value for 4 with α DF 0.05.

0.99997.

Area under curve = 1 \Rightarrow 1 - 0.99997

$$\Rightarrow \frac{0.00003}{2} = 0.000015.$$

This is p value = 0.000015.

p is 0.000015 $< \alpha$ 0.05 So we can reject H_0 .

2) In the population the Avg IQ is 100 with a SD of 15. Researchers want to test a new medication to see if there is +ve or -ve effect on intelligence or no effect at all. A sample of 30 participants who have taken the medication has a Mean of 140. Did the medication offer the intelligence or not with a C.I 95%?

4) H_0 : $\mu = 100$, $n = 30$, $\sigma = 15$, $\bar{x} = 140$. C.I = 95%
= 0.05

H_0 : The Avg ^{IQ of the} population is 100.

Since population SD (σ) is known - z-test.

H_A : ^{No} The Avg IQ of σ

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

$$z = \frac{140 - 100}{\frac{15}{\sqrt{30}}} = \frac{40}{\frac{15}{\sqrt{30}}} = \frac{40}{2.738} = 14.61.$$

At 95% Confidence (two-tailed test):

$$Z_{\text{critical}} = \pm 1.96.$$

$$\alpha = 1 - 0.95$$

$$= 0.05$$

$$Z_{\text{calculated}} = 14.61 > 1.96.$$

$$\frac{\alpha}{2} = \frac{0.05}{2} = 0.025$$

$$1 - 0.025 = 0.975$$

Reject H_0 .

$p < \alpha$ So we can reject

\therefore At 95% Confidence level, the Sample mean is

Significantly higher than the population mean. So, the

medication has a significant positive effect on

intelligence.

7/11/25

$$\mu = 100, \sigma = 15, n = 30, \bar{x} = 60, C.I. = 95\% = 0.05$$

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\frac{130 - 140}{140 - 130}}{\frac{15}{\sqrt{30}}} = -3.66.$$

$$= 0.00013.$$

$$p = 1 - 0.00013$$

$$= \frac{0.99987}{2}$$

$$= 0.499$$

$p > 0.05$ so accept the value.

1) Same above question the $\mu = 100$, $n = 30$, $\bar{x} = 140$, $s = 20$,
 $\alpha = 0.05$.

2) Here Sample SD is given so we have to do t-test.

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{140 - 100}{\frac{20}{\sqrt{30}}} = 10.95$$

$\text{DOF} = n - 1$
 $= 30 - 1$
 $\text{DOF} = 29$

$t >$ table value with $\text{DOF} = 29$ ~~accept H_0~~ & $\alpha = 0.05$
then reject H_0 else Accept H_0 .

$10.95 > 2.045$ so we can reject.

2) Credit Card Launch:

$n = 140$, $\bar{x} = 1990$, $\sigma = 2500$, $s = 2833$, $\text{C.I.} = 95\%$, $\alpha = 0.05$

We have to calculate the Interval Storage Range.

Default always the C.I. is 95%

$$\text{CI} = \bar{x} \pm z_{\alpha} \frac{\sigma}{\sqrt{n}} \quad \text{Lower Bound}$$

Upper Bound.

$$1990 - \frac{z_{0.05}}{2} \frac{2500}{\sqrt{140}}$$
$$= 1990 - z_{0.025} (211.29)$$

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

$$= 1990 - 414.1284$$

$$= 1576.$$

Upper Bound

$$\bar{x} + z_{\alpha} \frac{s}{\sqrt{n}}$$

$$= 1990 + \frac{z_{0.05}}{2} \left(\frac{2500}{\sqrt{140}} \right)$$

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

$$= 1990 + 414.12$$

$$= 2404.$$

∴ The Avg balance they are going to maintain after fullfledge launch is $[1576 \quad 2404]$

3) An^a quantitative 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% Confidence Interval about the mean?

a) $n = 25$, $\bar{x} = 520$, $s = 80$, $CI = 95\%$, $\alpha = 0.05$.

$$t = \bar{x} \pm t_{\alpha} \frac{s}{\sqrt{n}}$$

Lower Bound

$$\bar{x} - t_{\alpha} \frac{s}{\sqrt{n}}$$

$$= 520 - \frac{z_{0.05}}{2} \frac{80}{\sqrt{25}}$$

$$= 520 - (2.06) (16) = 520 - 31.36$$

Upper Bound

$$\bar{x} + t_{\alpha} \frac{s}{\sqrt{n}}$$

$$= 520 + \frac{z_{0.05}}{2} \frac{80}{\sqrt{25}}$$

$$= 520 + (2.06) (16)$$

$$= 520 + 31.36$$

$$= 552.$$

$$= 488.64$$

∴ the Avg balance they are going to maintain after fullfledged [487, 552].

Statistical tests:

