

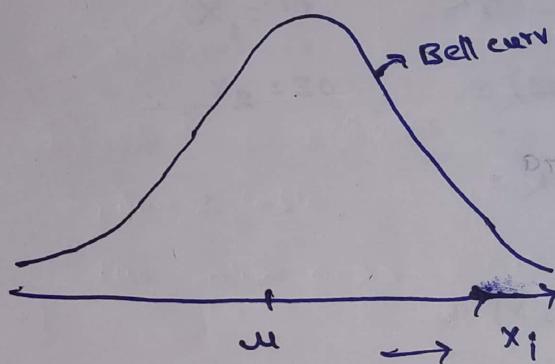
Central limit theorem: The central limit theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large number of samples taken from a population.

→ The central limit theorem says that the sampling distribution of the mean will always be normally distributed, as long as the sample size is large enough. Regardless of whether the population has a normal, poisson, binomial or any other distribution, the sampling distribution of the mean will be normal.

Let's say $x \sim N(\mu, \sigma)$

$$1 + d - d = \sigma$$

$$d = 1 + 1 - 1 = \text{sample size}$$



let's take few sample from population with size n

Here, $n = 1 \text{ or } 2 \text{ or } 3 \text{ or } 100 \dots$

$$S_1 = \{x_1, x_2, x_3, \dots, x_n\} = \bar{x}_1$$

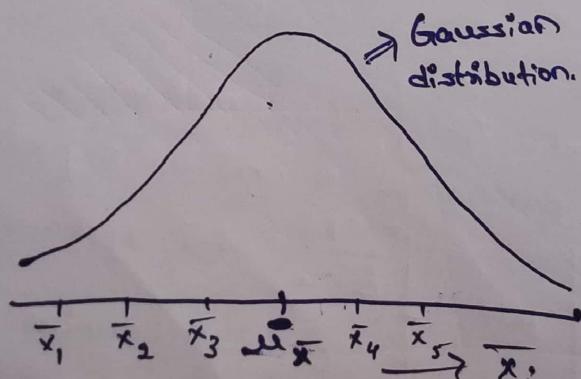
$$S_2 = \{x_1, x_2, x_3, \dots, x_n\} = \bar{x}_2$$

$$S_3 = \{x_1, x_2, x_3, x_4, \dots, x_n\} = \bar{x}_3$$

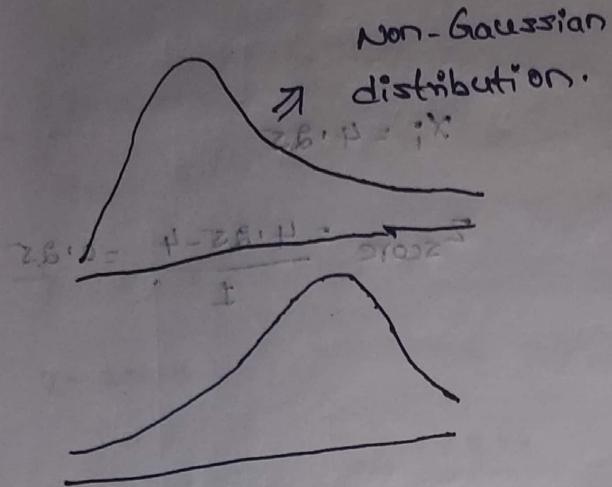
⋮

$$S_m = \{x_1, x_2, x_3, x_4, \dots, x_n\} = \bar{x}_m$$

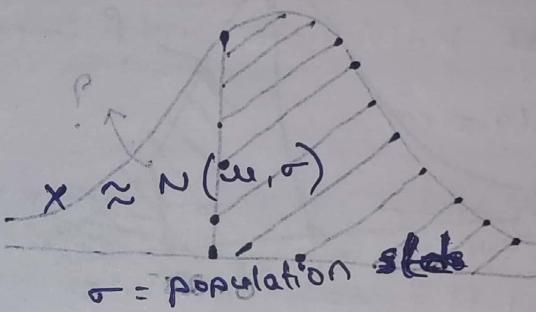
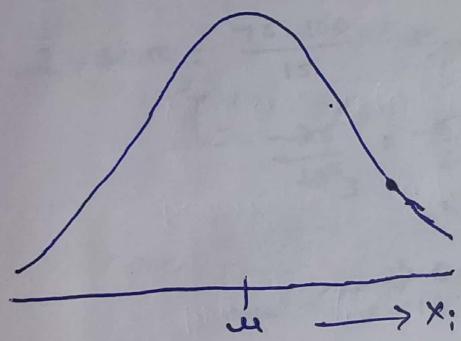
If we plot a graph with all the mean values of samples they also follow a Gaussian distribution.



$$\textcircled{2} \quad X \sim N(\mu, \sigma)$$



\rightarrow If X is not following Gaussian distribution we need to consider sample size ≥ 30 , in order to make sure that means of samples follow normal distribution.



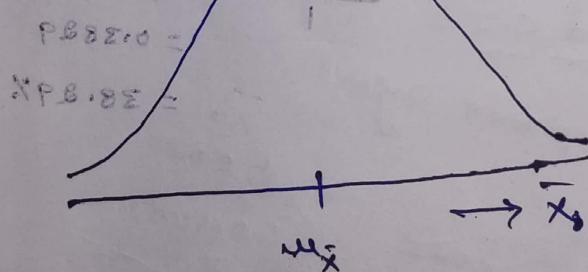
$\mu = population mean (approx)$

$1 + 8PZ \leq 0 \leftarrow Z = 0 \text{ value for}$

Sampling distribution of the mean

$$P(Z \leq 0) = 0.5 + \frac{1}{2} \Phi(0) = 0.5 + 0.5 = 0.5$$

$$P(Z \leq 0) = 0.5 + \frac{1}{2} \Phi(-1) = 0.5 - 0.5 = 0$$



$$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

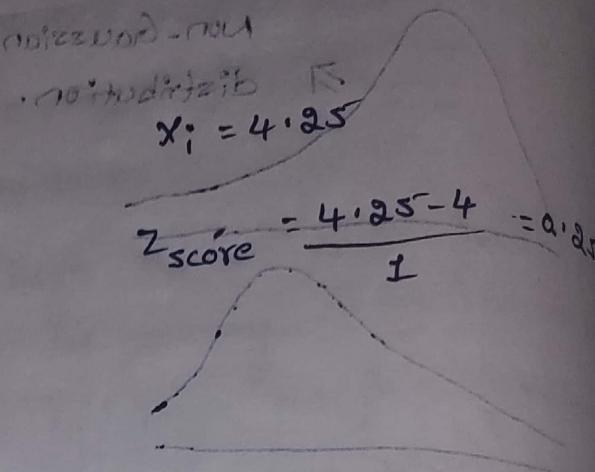
$\frac{\sigma}{\sqrt{n}}$ → standard error.
 $n = sample size$

$$\textcircled{1} \quad X \approx N(4, 1)$$

$$X = \{x_1, x_2, \dots, x_n\} \approx N(4, 1)$$

$$\bar{X} = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\} \approx N(4, \frac{1}{n})$$

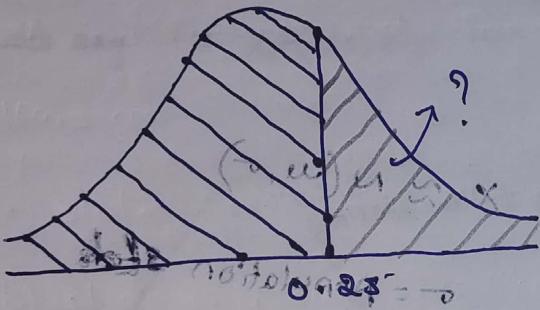
$$m\bar{X} = f(\bar{X}) \approx N(4, \frac{1}{n})$$



→ What percentage of scores falls above 4.25?

↳ How many standard deviations above the mean is 4.25?

↳ What percentage of scores falls above 0.25 standard deviations from the mean?

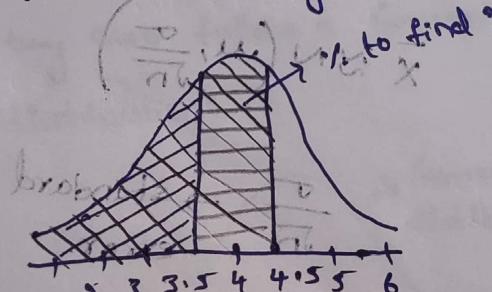


(longer) From Z-table
the value 0.25 → 0.59871

$$\text{Area a} \Rightarrow 1 - 0.59871 \\ = 0.4013$$

~~Area a = 40.13%~~

→ what percentage of scores lies between 3.5 to 4.5?



Size of sample = n

$$Z\text{-score} = \frac{4.5 - 4}{1} = +0.5 = 0.69147$$

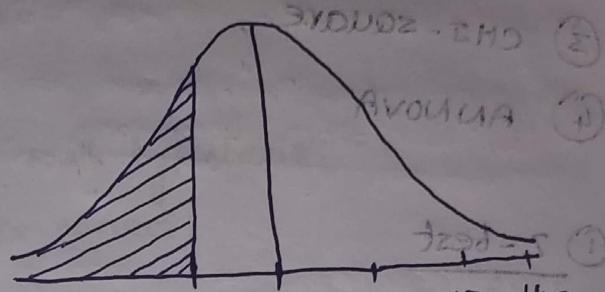
$$Z\text{-score} = \frac{3.5 - 4}{1} = -0.5 = 0.30857 \\ = 0.3829 \\ = 38.29\%$$



* In India the average IQ is 100, with a standard deviation of 15. What is the percentage of the population would you expect to have an IQ lower than 85?

$$\rightarrow \sigma = 15, \mu = 100$$

$$Z\text{-score} = \frac{85 - 100}{15} = -1 \\ \Rightarrow 0.1586$$



~~IQ < 85~~ $\Rightarrow 1 - 0.1586 = 0.8413$

i) $75 \leq IQ \leq 100$?

$$Z\text{-score} = \frac{75 - 100}{15} = -1.6666666666666667 \\ = -\frac{25}{15} = -\frac{5}{3} \\ = -1.6666666666666667$$

$$ZP(0) = 0.5 \Rightarrow 0.516$$

$$\text{Area} = 0.516 \times 2 = 1.032 \quad P(Z = 0) = 0.5$$

$$Z\text{-score} = \frac{100 - 100}{15} = 0$$

$$ZP(0) = 0.5 - 0.516 = 0.484$$

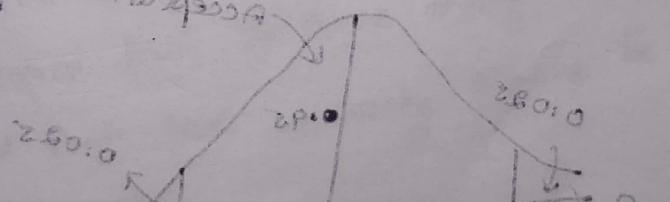
(ii) $z \geq z_0$ then $P(z \geq z_0)$

$z_0 = 1.645$

probability

$P(Z \geq 1.645) = 0.05$

highest probability



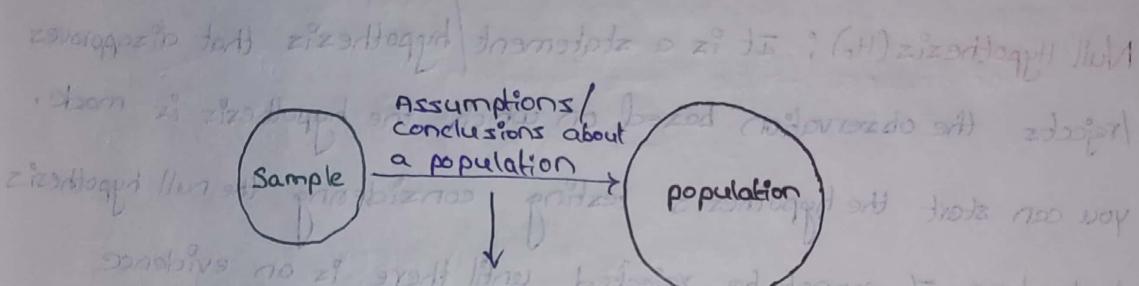
$$ZP(0) = 0.5 \\ ZP(1) = 0.5 + 0.3413 = 0.8413 \\ ZP(2) = 0.5 + 0.4772 = 0.9772 \\ ZP(3) = 0.5 + 0.4987 = 0.9987$$

Inferential Statistics

Statistical Inference

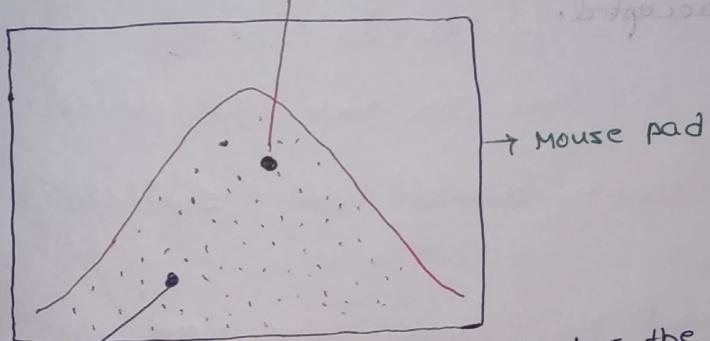
→ In statistical inference various sample statistics are analyzed to determine the population parameters.

Term	Denotes	Example
→ statistic	a quantity/property of sample	sample mean, sample variance
parameter	a quantity/property of population	population mean, population variance



To validate these assumptions, we use hypothesis testing.

p-value: It is the probability for null hypothesis to be true



$p=0.4 \Rightarrow$ out of all 100 touches, the no. of times is 40.

→ This p-value can determine probability.

Hypothesis Testing

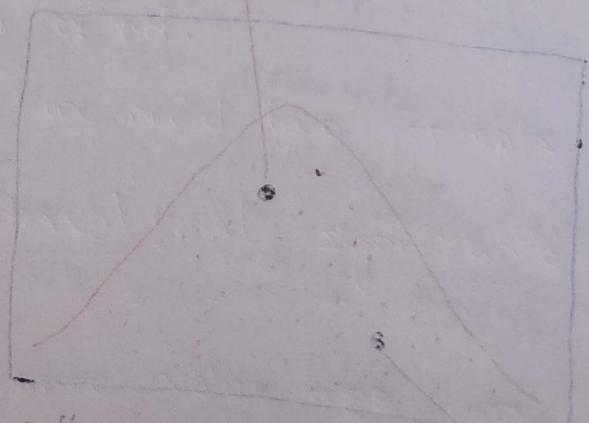
In statistics, hypothesis is a statement about the population which is represented by some numerical values.

Hypothesis testing deals with collecting enough evidence about the hypothesis. Then, based on evidence collected, the test either accepts or rejects null hypothesis.

→ each hypothesis test usually includes two competing hypotheses about the population. They are null hypothesis and alternate hypothesis.

Null Hypothesis (H_0): It is a statement/hypothesis that disapproves/rejects the observation based on which the hypothesis is made. You can start the hypothesis testing considering the null hypothesis to be true. It cannot be rejected until there is an evidence which suggests otherwise.

Alternate Hypothesis (H_1): It is a statement/hypothesis which is contradictory to the null hypothesis. If you find enough evidence to reject the null hypothesis, then the alternate hypothesis is accepted.



Steps involved in performing Hypothesis Testing

Step 1:-

figure out null hypothesis and Alternative hypothesis.

null hypothesis (H_0)

Alternate hypothesis (H_1)

significance level as

confidence interval

Step 2:- Decision boundary.

calculate the probability of getting the observed data (p value)

assuming null hypothesis to be true.

→ calculate the test statistic.

- 1) z-test
- 2) t-test
- 3) chi-square
- 4) ANOVA (F-test)

→ find out the problem is one tail or two tail.

→ find out confidence interval values

→ calculate the p-value.

Step 3:- conclusions.

① Fail to reject null hypothesis (when p value \geq significance level)

② Reject null hypothesis (when p value $<$ significance level)

Q) Test whether the coin is fair coin or not by performing 100 tosses.

① Null hypothesis - coin is fair $\rightarrow (H_0)$

② Alternate hypothesis - coin is not fair $\rightarrow (H_1)$

③ Experiment

④ conclusions

→ In the experiment performed we got following outcomes.

100 tosses \rightarrow 50 times head } coin is fair, both
50 times tails }

60 times head } coin is fair
40 times Tail }

30 times head } ?
70 times tail }

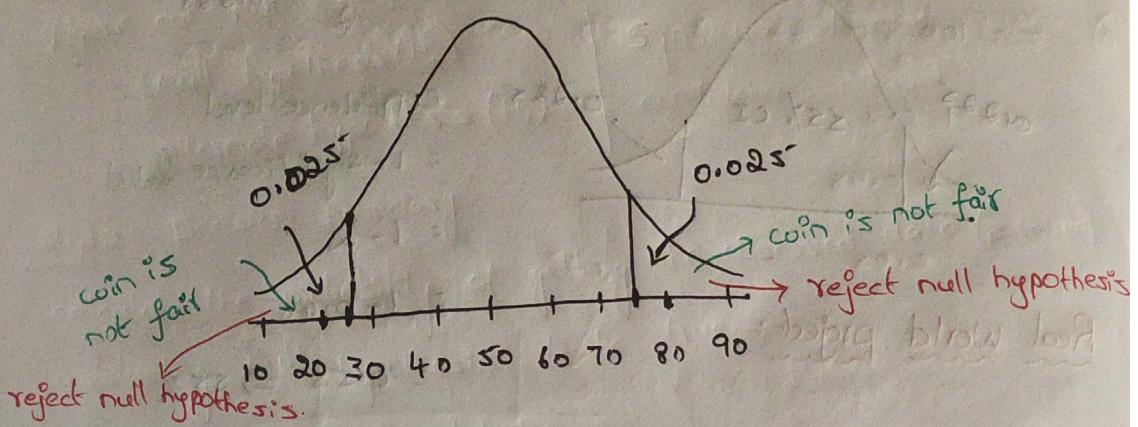
In this case we may get confusion to tell whether the coin is fair or not.

→ In this case we need two concepts one is "confidence interval" and "significance value".

→ Based significance value we can able to define confidence interval.

→ The significance value is given by domain expertise.

when I toss a coin 100 times, I may have chance to get heads 10 times head or 90 times head

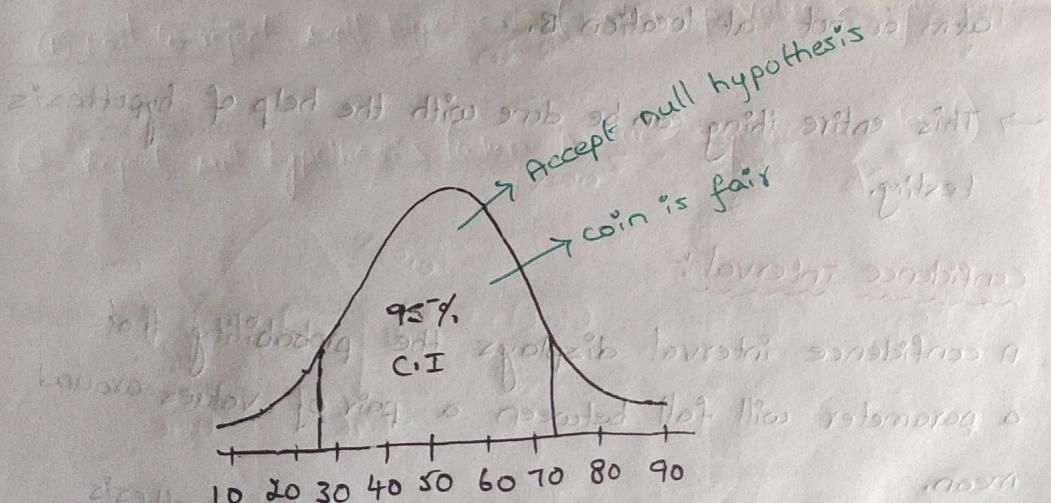


Conclusion:

Here 10, or 90 does not lies in range 20-80, so that null hypothesis is rejected.

⇒ coin is not fair.

If I perform a experiment I may have a chance to get 60 times head.



Conclusion:

Here 60 lies in range 20-80

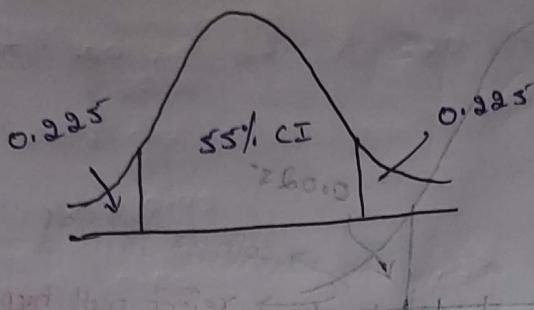
→ In this case we will accept null hypothesis

→ coin is fair.

for medical problem, $\alpha = 0.45$

if $\alpha = 0.45$, C.I. = ?

$$\frac{0.45}{2} = 0.225$$



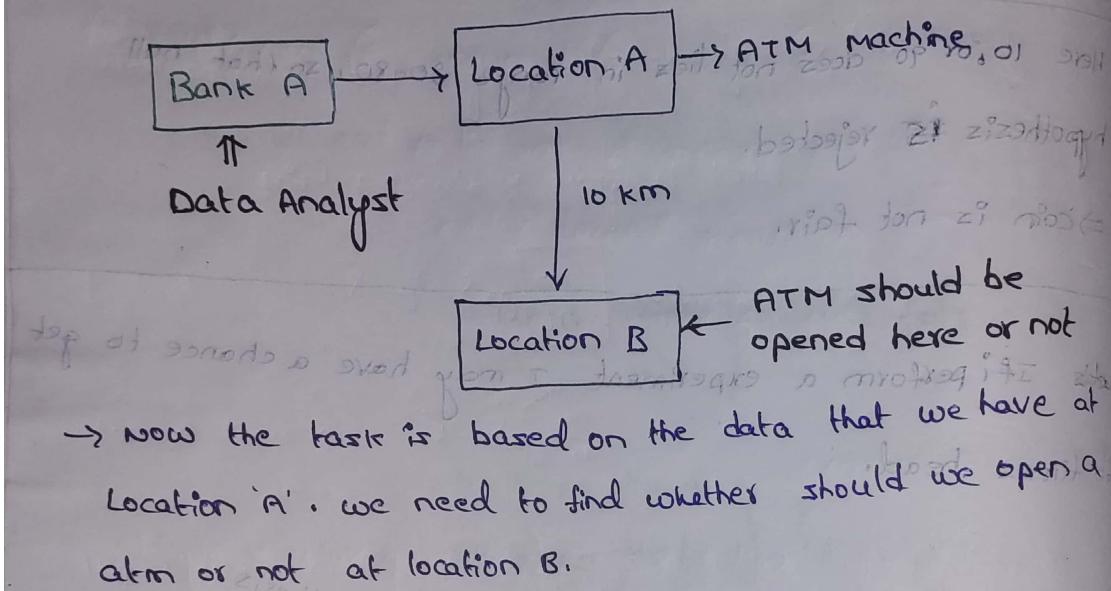
confidence level

$$= 1 - \text{significance level}$$

$$= 1 - \alpha$$

$$= 1 - 0.45$$

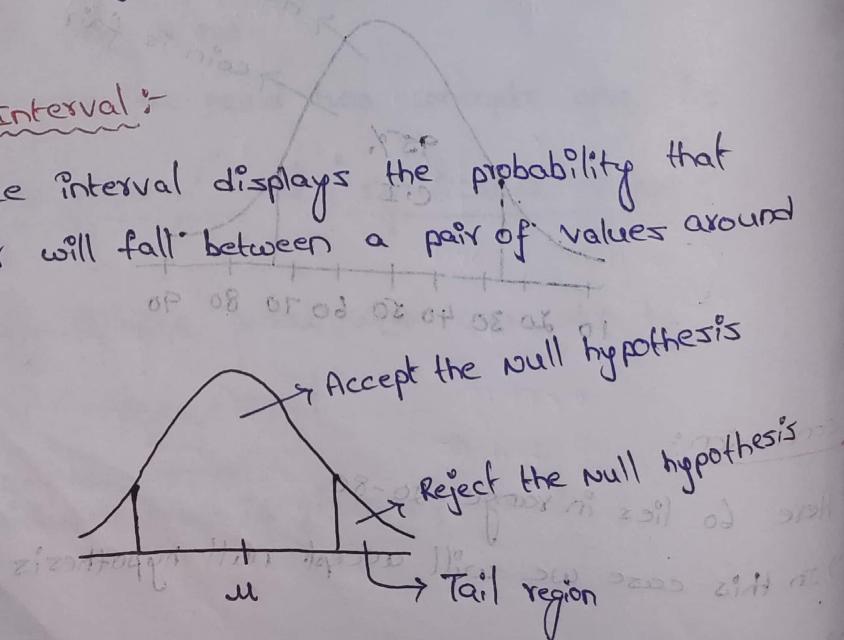
Real World project: $\alpha = 0.53 \Rightarrow 1 - \alpha = 0.47 \Rightarrow 47\%$



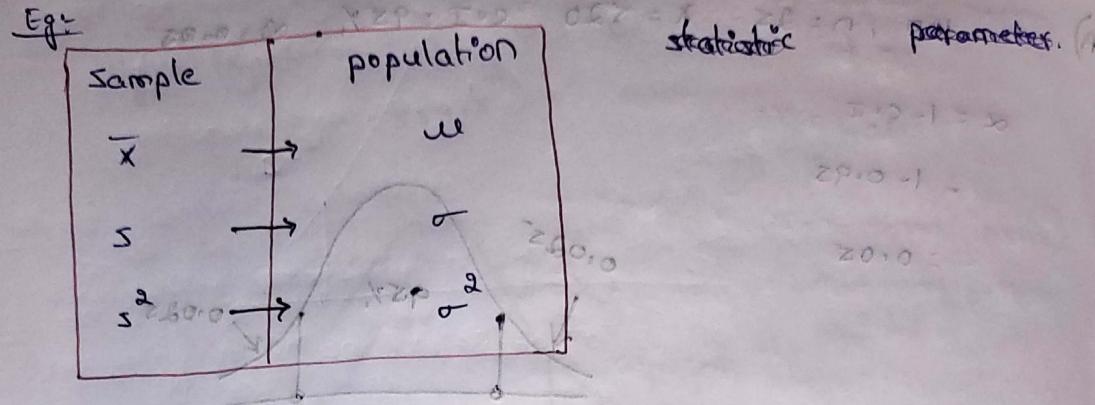
→ This entire thing can be done with the help of hypothesis testing.

Confidence Interval:

A confidence interval displays the probability that a parameter will fall between a pair of values around mean.



point estimate: the value of any statistic that estimates the value of a parameter is called point estimate.



sample mean \leftarrow from population mean \rightarrow \bar{x} μ

2.8, 2.9, 2.5, 3.5 \rightarrow 3

→ Sample (mean, variance, std deviation) are used to estimate population (mean, variance, std deviation).

Confidence Interval

point estimate \pm Margin of Error.

→ when population standard deviation is given:

$$C.I = \bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

M.O.E = $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$
 $z_{\alpha/2}$ = critical value of z distribution at significance level $\alpha/2$
 σ = population std deviation, n = no. of samples

C.I range for mean \bar{x}
 Lower C.I
 PE \pm MOE

to Higher fence C.I
 PE + MOE

C.I range
for mean

$$= \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \text{ to } \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

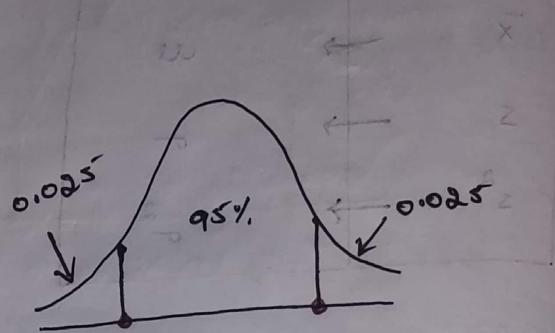
Q) On a quant test of CAT exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520. construct a 95% C.I about the mean?

A) $\sigma = 100, n = 25, \bar{x} = 520, C.I = 95\%, \alpha = 0.05$

$$\alpha = 1 - C.I$$

$$= 1 - 0.95$$

$$= 0.05$$



when population std is given $\{z\text{-score}\} \rightarrow z\text{-table}$

Point Estimate \pm Margin of Error.

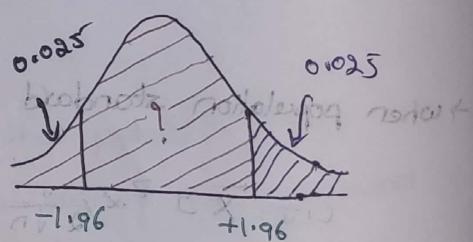
$$\text{Estimates of } \bar{x} \pm \frac{z_{\alpha/2} \sigma}{\sqrt{n}} \rightarrow \text{standard error.}$$

$$\text{Lower fence C.I} = \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$\text{Higher fence C.I} = \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Margin of error

$$z_{\alpha/2} = z_{0.025} \Rightarrow z_{0.0125}$$



In statistics $\frac{\sigma}{\sqrt{n}} = \text{std error}$
In probability σ/\sqrt{n} is called standard deviation
at 95% level significance
of error $\alpha = 0.05$, non-inclusive notation = 0.025

$$= 1 - 0.025$$

304 + 34

304 \pm 1.96 through z-table

$$= 0.9750 \text{ (from z-table)}$$

$$\frac{\sigma}{\sqrt{n}} \cdot z_{\alpha/2} + \bar{x} \text{ or } \frac{\sigma}{\sqrt{n}} \cdot z_{\alpha/2} - \bar{x} = \text{ Margin of error}$$

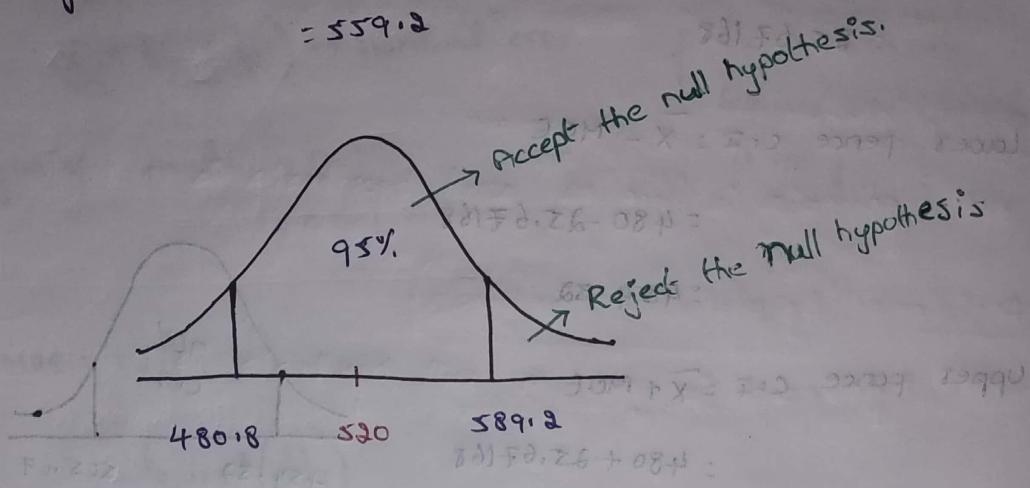
$$\text{lower fence} = 520 - 1.96 \times \frac{100}{\sqrt{25}}$$

$$= 520 - 1.96 \times 20$$

$$= 480.8$$

$$\text{Higher fence} = 520 + 1.96 \times 20$$

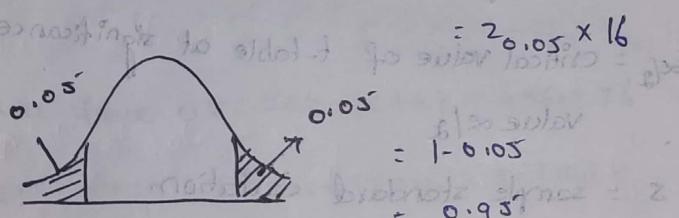
$$= 559.2$$



- Q) $\sigma = 80, n = 25, CI = 0.90, \bar{x} = 480$. construct a 90% CI about mean. $\alpha = 1 - CI = 1 - 0.9 = 0.1$

$$MOE = Z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}} = \frac{Z_{0.05}}{2} \times \frac{80}{\sqrt{25}}$$

Margin of error = \bar{x} marks



0.95 lies between 0.9450 and 0.9505

i.e. $\Rightarrow 1.604$ and 1.605

x	0.9450	0.95	0.9505	0.95	0.9505
y	1.604	2.05	1.605	2.05	1.605

$$Y - Y_1 = \frac{Y_2 - Y_1}{X_2 - X_1} (X - X_1) \rightarrow \text{Interpolation}$$

or $= 2.05$ and 2.05 , out of $6 = 1.6 + 1.6 \times 0.05$

$$Z_{0.05} = 1.604 + \frac{1.605 - 1.604}{0.95053 - 0.94950} \times (0.95 - 0.94950)$$

$$Z_{0.05} = 1.60448$$

$$MOE = 1.60448 \times 16$$

~~= 25.68~~

$$= 25.67168$$

$$\text{Lower fence } CI = \bar{x} - MOE$$

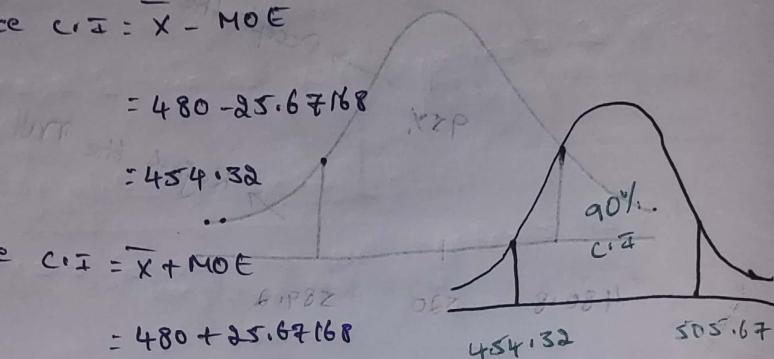
$$= 480 - 25.67168$$

$$= 454.32$$

$$\text{Upper fence } CI = \bar{x} + MOE$$

$$= 480 + 25.67168$$

$$= 505.67$$



t-test

$$CI = \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}}$$

where \bar{x} = mean of sample

$t_{\alpha/2}$ = critical value of t-table at significance

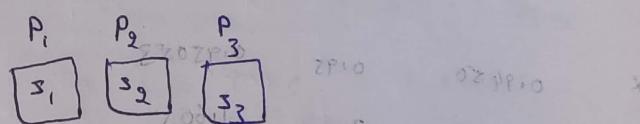
value $\alpha/2$

s = sample standard deviation.

n = no. of samples

Degree of freedom (d.f.) = $n-1$

suppose 3 seats are available



3 person

- 1st person (P_1) → 3 choices to sit
- 2nd person (P_2) → 2 choices to sit
- 3rd person (P_3) → no choice must sit on last seat

$d.f. = n-1 = 3-1 = 2 \Rightarrow$ only two people have choice to choose seat

→ On the Quant test of cat exam, a sample of 25 test takers has a mean of 520 with a sample standard deviation of 80.

construct 95% C.I about the mean?

Ans: $\bar{x} = 520, s = 80, \alpha = 0.05, n = 25, C.I = 0.95$
when sample standard deviation is given we have to use t-test=t-table

$$\bar{x} \pm t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right) \rightarrow \text{standard error. } \quad \left\{ \begin{array}{l} \text{t-test we are using} \\ \text{because population std} \\ \text{is not given} \end{array} \right.$$

$$\text{Degree of freedom} = n - 1 \\ = 25 - 1 \Rightarrow 24$$

$$MOE = t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right)$$

$$= t_{0.025} \times \left(\frac{80}{\sqrt{25}} \right)$$

(From t-table) $t_{0.025}$ and dof = 24

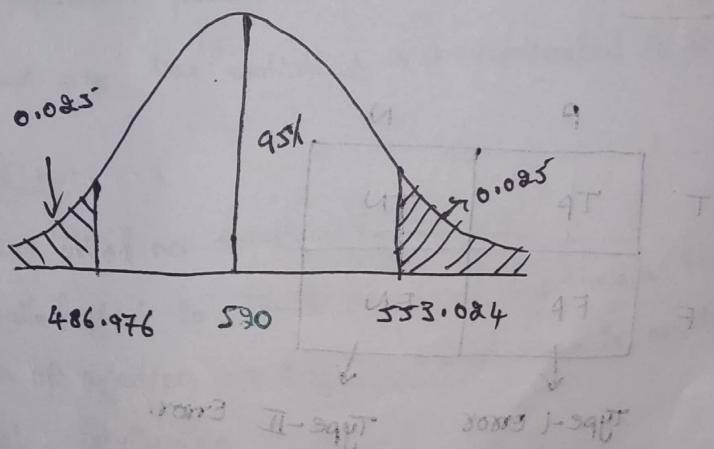
$$MOE = t_{0.025} = 2.064$$

$$MOE = 2.064 \times 16$$

$$= 33.024$$

lower fence C.I = $520 - 33.024 = 486.976$

upper fence C.I = $520 + 33.024 = 553.024$



Type-I and Type-II Error

Null hypothesis (H_0) = coin is fair

Alternate hypothesis (H_1) = coin is not fair.

Reality check:

null hypothesis is true or null hypothesis is false.

Decision's {By performing Experiments}

null hypothesis is true or null hypothesis is false.

outcome : 1

we reject the null hypothesis, when in reality it is false \rightarrow yes (Good decision)

outcome : 2

we reject the null hypothesis when in reality is true \rightarrow Type-I Error

outcome : 3

we accept the null hypothesis when in reality it is false \rightarrow Type-II Error

outcome : 4

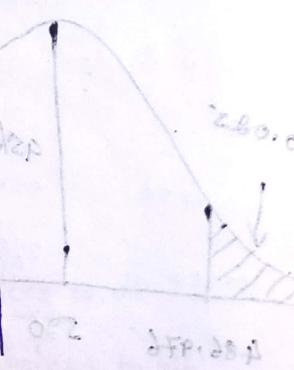
we accept the null hypothesis when in reality it is true \rightarrow yes (Good decision)

Confusion Matrix:

	P	N
T	TP	FP
F	FN	TN

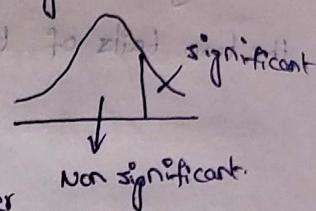
Type-I Error

Type-II Error.



one-tailed test

- It is based on uni-directional hypothesis.
- one tailed test is a hypothesis test in which the region of rejection appears on one end of sampling distribution.
- critical values are boundary between non significant and significant results in a hypothesis testing.



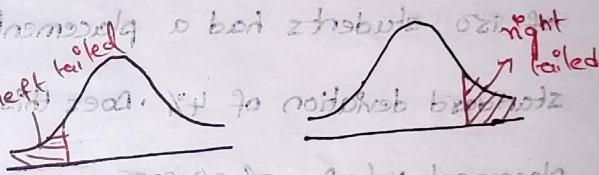
→ one-tailed test represents estimated parameter

is greater or less than the critical value.

→ when the sample tested falls in the region of rejection either left or right side, at this time it leads to acceptance of alternate hypothesis rather than null hypothesis.

In this statistical hypothesis test, all the critical region related to is placed any one of the two tails.

one-tailed test can be:



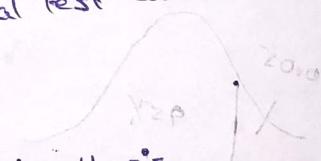
Left-tailed test:

when the population parameter is believed to be lower than the assumed one, the hypothesis test carried out is the left tailed test.



Right-tailed test:

when the population parameter is supposed to be greater than the assumed one, the statistical test conducted is a right tailed test.

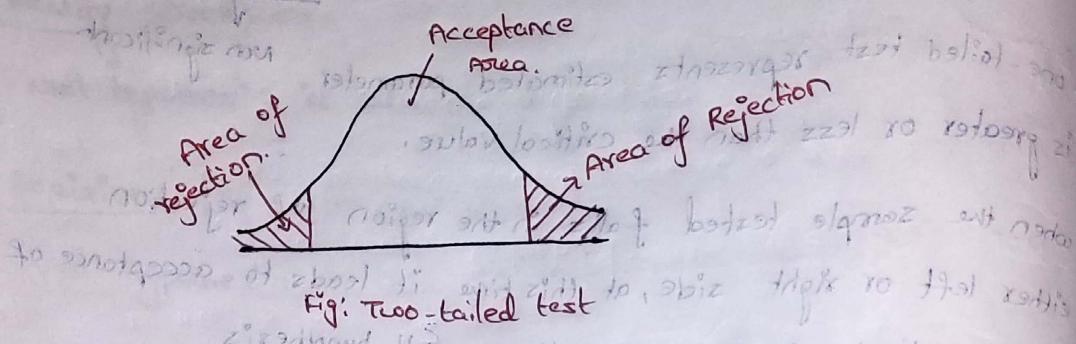


Two tailed test:

It is also called non-directional hypothesis.

The two tailed test is described as a hypothesis test, in which the region of rejection or say the critical area is on both ends of the normal distribution.

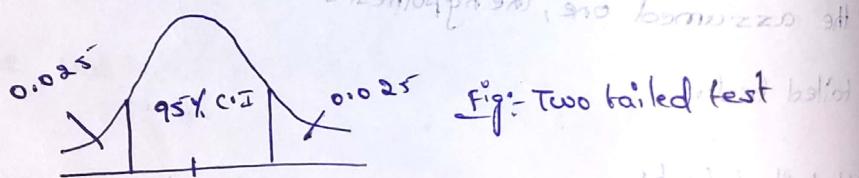
- It determines whether the sample tested falls within or outside of a certain range of values.
- Therefore the alternate hypothesis is accepted in place of the null hypothesis, if the calculated values falls in either of the two tails of the probability distribution.



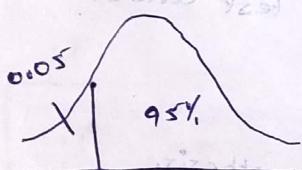
- Q) colleges in karnataka has an 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88% with a standard deviation of 4%. Does this college has different placement rate?

$$\text{standard deviation of sample} = \sigma = 0.005$$

- The true placement rate is may be less than or greater than 85%.



- If placement is less than 85%,



Difference between one tailed and two tailed test.

one-tailed test	two-tailed test.
→ used to test a directional hypothesis.	→ used to test a non-directional hypothesis.
→ The null hypothesis is rejected if the test statistic falls entirely on one side of sampling distribution.	→ The null hypothesis is rejected if the test statistic falls in either tail of the sampling distribution.
→ The critical value or p-value is only calculated for one tail of the distribution.	→ The critical value or p-value is calculated for both tails of the distribution.
→ we use either " $<$ " or " $>$ " sign for (H_1).	→ we use \neq sign for H_1 .
→ If H_1 specifies any direction then we use one tailed hypothesis.	→ if no direction given then we use two tailed hypothesis.

When to use z-test?

- when sample size ($n \geq 30$) and population standard deviation is given.
- (1) when sample standard deviation or sample standard deviation is given.
- (2) either of them should be given.

When to use t-test?

- when sample size $n < 30$
- (1) when sample size of sample means given is $n = 10$ to 20 .
- (2) when sample standard deviation is given.

$$Z_{\text{cal}} = Z_{\text{crit}} = \bar{x} - \mu_0 / \sigma / \sqrt{n}$$

Q: In a hypothesis test

+ 2 marks

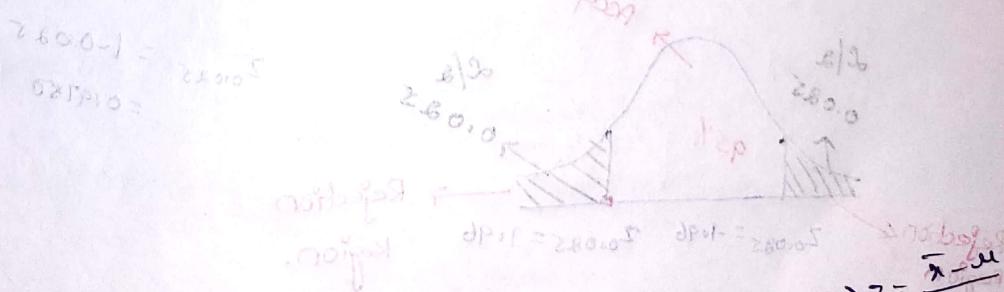
$H_0: \mu = (\text{H})$ zizzitoggi lura

$H_A: \mu \neq (\text{H})$ zizzitoggi laram

3 marks

$$\alpha = 0.05 \Rightarrow Z_{\text{crit}} = 1.96$$

prohibited noise level



$$\text{when } n=1 \Rightarrow Z = \frac{\bar{x}-\mu}{\sigma}$$

$$\text{when } n=30 \Rightarrow Z = \frac{\bar{x}-\mu}{\sigma/\sqrt{n}}$$

using sd blanda cariuska kahita sa gano po bts matalugay