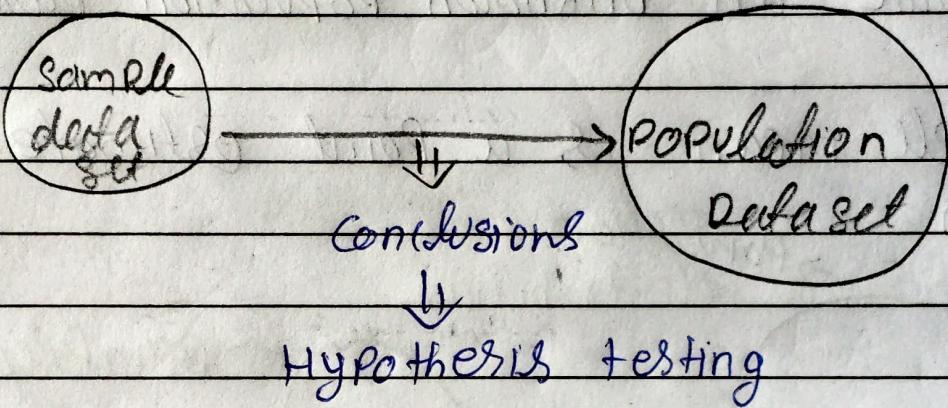


Inferential Statistics

Hypothesis Testing and Its mechanism:

Inferential Stats: \Rightarrow some conclusions or inferences.



Hypothesis testing mechanism:

(i) Null Hypothesis (H_0):

The assumptions you are beginning with.

(ii) Alternate Hypothesis (H_1):

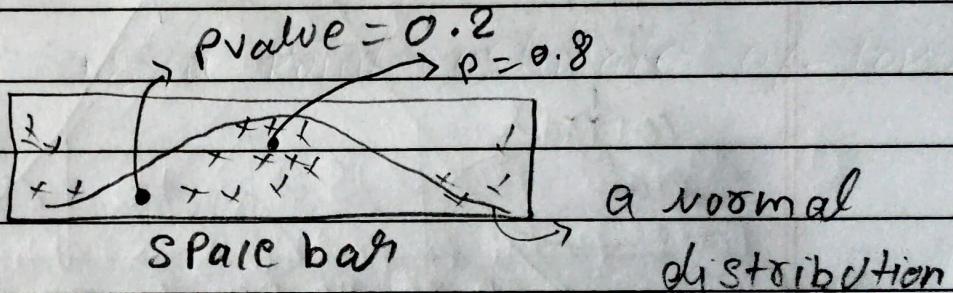
The opposite of null - Hypothesis.

(iii) Experiments: \rightarrow Statistical Analysis. \rightarrow collect the proofs.

(iv) Accept the Hypothesis or reject the hypothesis.

P-Value:

The P-Value is a number, calculated from a statistical test, that describes how likely you are to have hypothesis testing to help decide whether to reject the null hypothesis.



so here $p\text{value} = 0.2$ means if we touch the Space Bar 100 times, out of 100 touches we touch ~~more~~ around 20 times in that region.

Eg:- I am tossing a coin 100 times.

if $P(H) = 0.5, P(T) = 0.5 \Rightarrow$ good

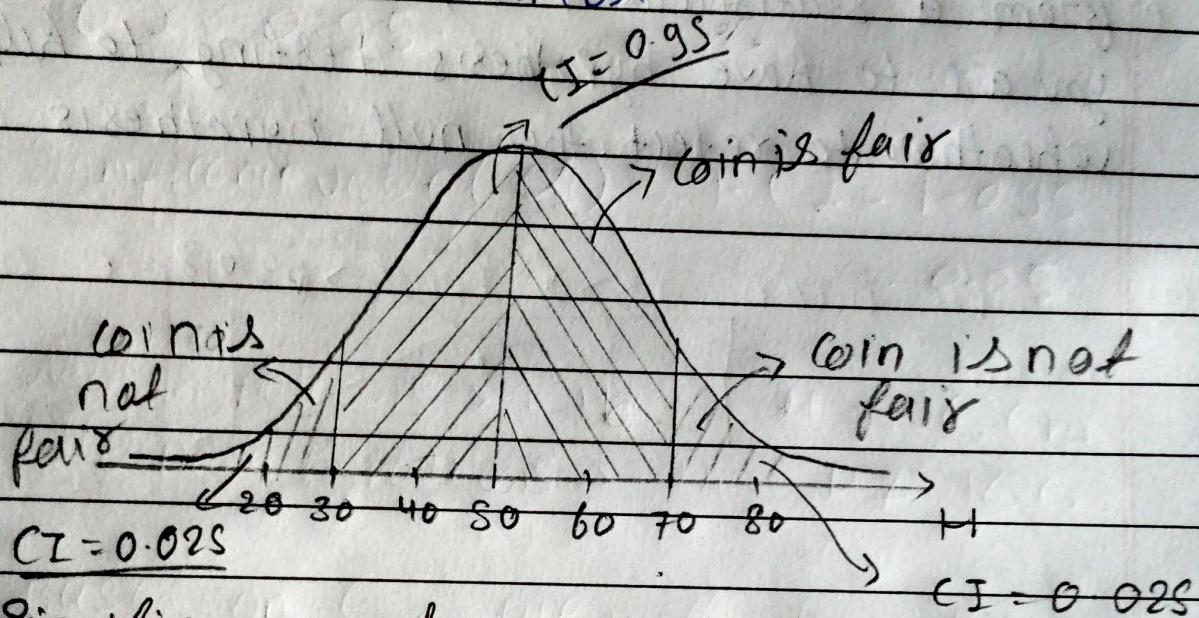
$\Rightarrow P(H) = 0.6, P(T) = 0.4 \Rightarrow$ good

$\Rightarrow P(H) = 0.7, P(T) = 0.3 \Rightarrow$ something is wrong.

① Null hypothesis:- coin is fair.

② Alternate hypothesis:- coin is not fair

③ Experiments: 100 times.



④ Significance value:

$$\alpha = 0.05$$

\rightarrow (area)

$$(C) \text{ confidence Interval} = 1 - 0.05 = 0.95$$

$P < \text{Significance value} \Rightarrow \text{Reject the null}$

\rightarrow hypothesis

Type of tests:

- ① Z-test = } \Rightarrow for Average data.
- ② t-test } \Rightarrow for Average data.
- ③ CHI-SQUARE test \Rightarrow for categorical data.
- ④ ANNOVA test \Rightarrow variance related data/analysis

Z-test:-

To apply these thing or test we have
2 conditions:

① Sample size > 30

② We know about the Population Std.

Q. The avg heights of all residents in a city is 168 cm with a ($\sigma = 3.9$) Std. A doctor believes the mean to be different. He measures the height of 36 individuals and found the avg height to be 169.3 cm

- ① State Null and Alternative hypothesis
- ② At a 95% confidence level, is there enough evidence to reject the null hypothesis.

80%

Soln

$$\mu = 168 \text{ cm}$$

$$\sigma = 3.9$$

$$n = 36$$

$$\bar{x} = 169.5 \text{ cm}$$

$$C.I = 0.95$$

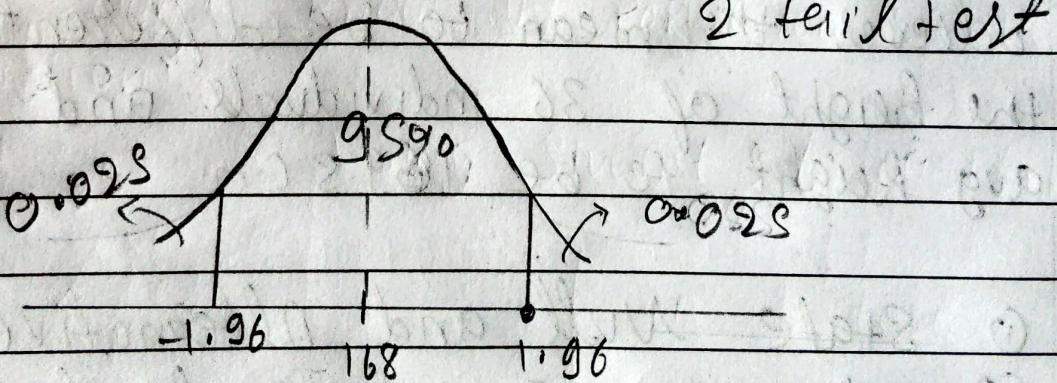
$$\text{Significance value } (\alpha) = 1 - C.I = 1 - 0.95$$

$$\alpha = 0.05$$

① Null hypothesis $\Rightarrow H_0: \mu = 168 \text{ cm}$

Alternate hypothesis $\Rightarrow H_1: \mu \neq 168 \text{ cm}$

\Rightarrow Based on C.I, we will draw the decision boundary.



$$Z\text{-Score} \Rightarrow 1 - 0.025 = 0.9725$$

So based on the given area we get the +1.96 Z-Score

If Z is less than -1.96 or greater than 1.96 , Reject the null hypothesis.

Z-test \Rightarrow

$$Z\text{-Score} = \frac{\bar{X} - \mu}{\sigma}$$

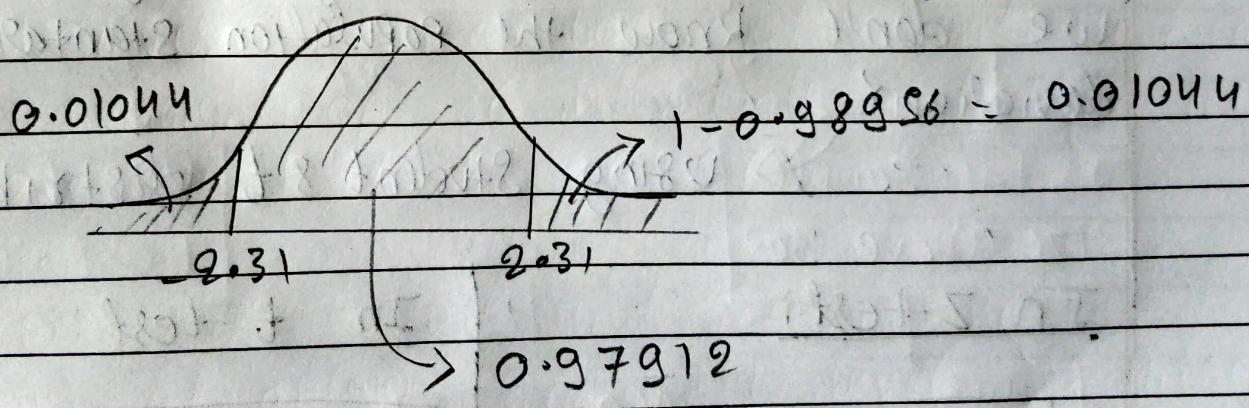
$$Z_d = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

$$Z_d = \frac{169.5 - 168}{3.9 / \sqrt{36}}$$

} apply central limit theorem
 $\sigma \Rightarrow \frac{\sigma}{\sqrt{n}}$

$$= \frac{1.5}{0.65} = 2.31$$

$2.31 > 1.96$ Reject the null hypothesis.



$$(0.98756 - 0.01044)$$

$$\text{P-value} = 0.01044 + 0.01044 \Rightarrow 0.02088$$

$$P < 0.05$$

$0.09088 < 0.05 \Rightarrow$ Reject the Null Hypothesis

\Rightarrow final conclusion the Avg Heights $\neq 168\text{cm}$

The avg height seems to increasing based on sample height.

Student t-test:

In z-test when we perform any analysis using z-score we require σ (population standard deviation) \rightarrow is already known.

\Rightarrow How do we perform any analysis when we don't know the population standard deviation?
 \Rightarrow using Student's t distribution

In z-test:

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

Population Std

In t-test:

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

Sample Std

In t-test we use a parameter
Degree of freedom.

$$\left\{ \text{Dof} = n - 1 \right\}$$

n = no of useful sample

$$\alpha = CI = 95\%$$

$$\alpha = 0.05$$

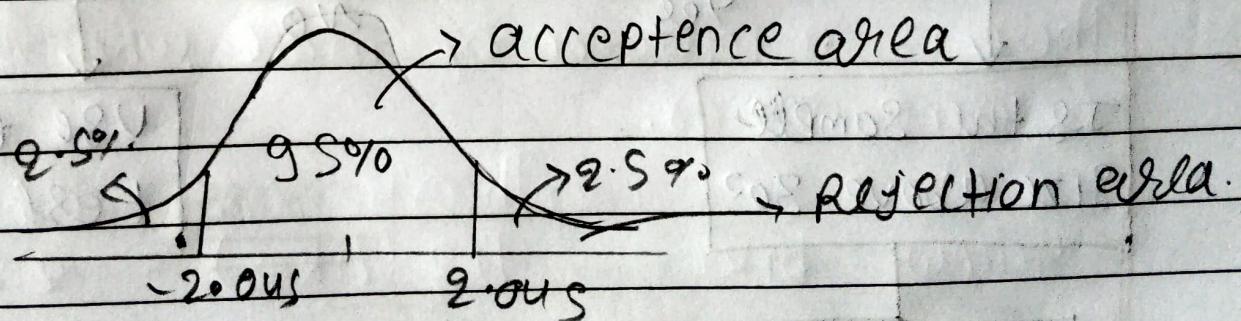
$$\mu = 100, S = 20$$

$$n = 30, \bar{x} = 140$$

Null hypothesis :- $H_0 = \mu = 100$

Alternate hypothesis : $H_1 = \mu \neq 100$

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



If my t-test is less than -2.045 or greater than 2.045, we have to reject the null hypothesis.

t-test

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

$$t = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65} = 10.96$$

Since $t = 10.96 > 2.045$ we have to reject the null hypothesis.

when to use z-test and t-test

Do you know the population std(σ)?

↓
Is the sample size above 30?

Yes

No

↓
use t-test

↓
z-test

↓
No

↓
t-test

Bayes' theorem:

Bayesian statistics is an approach to data analysis and the parameter estimation based on Bayes' theorem.

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

$$\Rightarrow P(A) * P(B|A) = P(B) * P(A|B)$$

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

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

Bayes theorem

Dataset

Size of house, No of rooms

location

x_1

independent

↑

x_2

x_3

dependent

price

of P

y

$$P(y/x_1, x_2, x_3) = \frac{P(y) * P(x_1, x_2, x_3/y)}{P(x_1, x_2, x_3)}$$

→ Bayes theorem on your given dataset

CHI SQUARE TEST:-

The Chi Square test for goodness of fit test claims about population proportion.

It ~~is~~ is a non-parametric test that is performed on categorical data [ordinal and nominal] data.

eg:- There is a population of mall who likes different - colors.

	Theory	Sample
yellow bike	113 (Likes)	22

Red bike	113 (Likes)	17
----------	-------------	----

orange bike	113 (Likes)	59
-------------	-------------	----

theory categorical distribution observed categorical distribution

⇒ Goodness of fit test we are trying to implement

#

$$\chi^2 = \frac{(O - E)^2}{E}$$

$$\chi^2 = \frac{(O - E)^2}{E}$$

Analysis of variance (ANOVA) :-

⇒ ANOVA is a statistical method used to compare the means of 2 or more groups.

i) Factors (variable)

ii) Levels

e.g.: medicine (factor)

Dosage { 10 mg 5 mg 15 mg } ⇒ Levels

mode of Payment (factor)

Levels \leftarrow [GPay Phone Pay IMPs]

Assumptions in ANOVA:-

- ① Normality of sampling distribution of mean \rightarrow
The distribution of sample mean is normally distributed
- ② Absence of outliers \rightarrow Outlying score need to be removed from the data set.
- ③ Homogeneity of variance \rightarrow population variance in different levels of each independent variables are equal.
$$\left[\sigma_1^2 = \sigma_2^2 = \sigma_3^2 \right]$$
- ④ Samples are random and independent.

Hypothesis testing in ANOVA:-

① Null Hypothesis: - $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_K$

② Alternate hypothesis: - $H_1:$ At least one of the sample mean is not equal

F-test:-

$$\left\{ F = \frac{\text{Variance between samples}}{\text{variance within samples}} \right\}$$

[Variance between samples]		
X ₁	X ₂	X ₃
1	6	5
2	7	6
4	3	3
5	2	2
3	1	4
$\bar{x}_1 = 3$		$\bar{x}_2 = 19/5$
$\bar{x}_3 = 4$		

like the data is spread with respect to mean.

Q.

	15 mg	30 mg	45 mg	
9	7	4		
8	6	3		
7	6	2		
8	7	3		
8	8	4		
9	7	3		
8	6	2		

12 hours random

① Define Null and alternate hypothesis.

$$H_0: \mu_{15\text{mg}} = \mu_{30\text{mg}} = \mu_{45\text{mg}}$$

$$H_1: \text{not all } \mu \text{ are equal}$$

② Significance score: $\alpha = 0.05$

$$C.I = 0.95$$

③ calculate the Degree of freedom

$$N = 21, \alpha = 3, n = 7$$

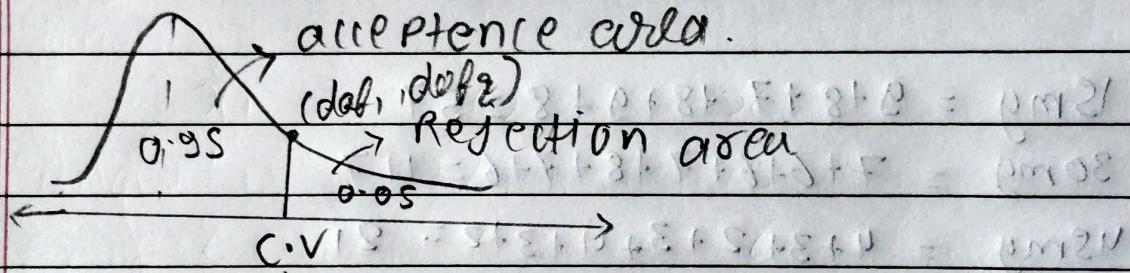
$$\text{df (between)} = \alpha - 1 = 3 - 1 = 2$$

$$\text{def (within)} = N - \alpha = 21 - 3 = 18$$

$$\text{def (total)} = N - 1 = 21 - 1 = 20$$

$(2, 18)$ } if $\alpha = 0.05$ it is a critical value.
 \downarrow \downarrow def_2 \downarrow (c.v)
 def_1

④ Decision boundary:



↳ based on the F-table

Value of c.v at $(2, 18) = 3.5546$

Null:

If F is greater than 3.5546 , reject the Null Hypothesis.

⑤ Calculate the F-test:

	sum of square	df	mean square	F
→ between	98.67	2	49.34	
→ within	10.29	18	0.54	
→ total	108.96	20	5.44	

$$\textcircled{1} \quad SS_{\text{between}} = \frac{\sum (\sum a_i)^2}{n} - \frac{T^2}{N}$$

$$15mg = 9+8+7+8+9+8 = 57$$

$$30mg = 7+6+6+7+8+7+6 = 47$$

$$45mg = 4+3+2+3+4+3+2 = 21$$

$$\Rightarrow \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57^2 + 47^2 + 21^2]}{21}$$

$$\Rightarrow 98.67$$

$$\textcircled{2} \quad SS_{\text{within}} \Rightarrow \sum y^2 - \frac{\sum (\sum a_i)^2}{n}$$

$\sum y^2$ = sum of squares of all the values.

$$\sum y^2 = 9^2 + 8^2 + 7^2 + \dots \\ = 853$$

$$\Rightarrow 853 - \frac{[S7^2 + 47^2 + 21^2]}{7}$$

$$\Rightarrow 10.29$$

\Rightarrow mean square = $\frac{SS}{dof}$

$$F\text{-test} = \frac{MS(\text{between})}{MS(\text{within})}$$

$$\Rightarrow \frac{49.34}{0.84}$$

$$\Rightarrow \underline{\underline{86.56}}$$

$\Rightarrow 86.56 > 3.5546$ so we have to reject the null hypothesis.