

→ Inferential Statistics Date :- 30/3/22
~~~~~ \* ~~~ 11:00 AM

## \* Central Limit Theorem :- For Continuous Data

Example :-

Temperature @ Deshaipet (My home) @ 6 pm

| Day | Jan 1 | 27.  | Avg = 30.5 |
|-----|-------|------|------------|
| *   | Jan 2 | 29   |            |
| *   | Jan 3 | 34   |            |
| *   | Jan 4 | 30   |            |
| *   | Jan 5 | 31.5 |            |
| *   | Jan 6 | 32   |            |
| *   | Jan 7 | ?    |            |

(What is Jan 7 Temperature)

$$\text{avg}(\bar{x}) = 30.5$$

\* What probability of 30.5 value for Jan 7

$$P(30.5) = 0, \leq = 1$$

$$P(30.5) = \frac{1}{\text{infinite}}$$

$$\approx 0, " "$$

or Single Value probability is "Zero".  
So, we adding Error Both sides,

We have "add" Error

$$\bar{x} \pm \text{std.error} \rightarrow \text{standard.error}$$

$$\bar{x} = (30.5 \pm 5)$$

Center value

Two side error

$$[\bar{x} - \text{std.error}, \bar{x} + \text{std.error}]$$

$$\therefore \text{std.error} = \frac{\sigma}{\sqrt{n}}$$

$$\bar{x} + \frac{\sigma}{\sqrt{n}}$$

Based on  
sample  
collection

Central Limit theorem

$$[\bar{x} - \frac{\sigma}{\sqrt{n}}, \bar{x} + \frac{\sigma}{\sqrt{n}}]$$

Mean

Sample

Predicting  
with in  
Range

\* For Continuous Data, In order to predict, The

Range, we apply Central Limit Theorem.

Not For "Discrete Variable"

Because Discrete data doesn't

have standard deviation.

with what confidence, (or) what probability ?

Confidence Interval:

$$\bar{X} - Z_{1-\alpha} \frac{\sigma}{\sqrt{n}}, \quad \bar{X} + Z_{1-\alpha} \frac{\sigma}{\sqrt{n}}$$

\*  $Z_{1-\alpha} = Z_{\text{confidence}}$

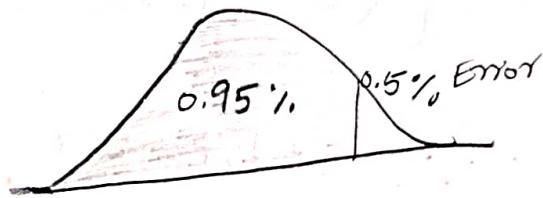
$$= Z_{0.95} \Rightarrow 5\% \text{ Error}$$

$$= Z_{0.90} \Rightarrow 10\% \text{ error}$$

From

$Z$  Table,  $Z_{95\%} = 1.65$

|      |      |      |      |       |       |              |
|------|------|------|------|-------|-------|--------------|
| $Z$  | 0    | 0.01 | 0.02 | 0.03  | 0.04  | 0.05         |
| +1.6 | .945 | .946 | .947 | .9484 | .9495 | .9505<br>95% |



with 95% confidence

$$\therefore \left[ \bar{X} - 1.65 \frac{\sigma}{\sqrt{n}}, \quad \bar{X} + 1.65 \frac{\sigma}{\sqrt{n}} \right]$$

What is  $(1 - \alpha)$  in Confidence Interval?

A:-  $\alpha = \text{Error}$

$1 - \alpha = \text{Confidence}$

Hypothesis testing → Statistical Testing

We make conclusions ⇒ population data based on sample data.

\* Continuous Variable :-

→ 1 Sample Test → 1. Sample Z Test  
1. Sample T Test

→ 2 Sample Test → 2. Sample Test For Equal Variance  
2. Sample Test For Unequal Variance.

→ Anova Test

\* Discrete Variable :-

1 proportion Test

2 proportion Test

Chi-square Test

Example:- "SRK" Teaching Feed Back = ~~5~~ rating

⇒ 15 students (response) - Sample Data

⇒ Avg. rating = 4.9

⇒ population Data = 50 students  
→ maybe  
1.0  
2.0  
2.5 ratings in population

To check by

(Statistical Testing) Hypothesis Testing

(Conclusions / inferences) ⇒ Based on Sample ( $\bar{x}$ ),  $\sigma$ ,  $s$   
Pop. Avg, Pop. Var, P. op. std.

+ (Estimating population) ( $M$ )

For Sample, Consider same for population data

(Sample ki Edi aye thay, apply chesthamo, Same population  
apply cheyali, Kuda aye, Apply cheyali)

Inferential Statistics

→ Your Inferencing Game one,  
Based on something.

conducted on "Samples".

→ Statement

→ Hypothesis testing :-

- \* it is used To determine whether a statement about "value of a population" parameter **should** (or) **shouldn't** be **rejected**.

Example :-

In Indian Court of Law,

- \* person is innocent
- \* person is criminal

In what criteria, we can Judge him.

So, we have.

\* Basic Assumption  $\Rightarrow$  Null hypothesis  $[H_0]$

\* Evidence of crime  $\Rightarrow$  Alternative hypothesis  
 $[H_1 / H_a]$

→ Hypothesis Testing

• Null Hypothesis  $[H_0]$  = Person is innocent

• Alternate Hypothesis  $[H_1 / H_a]$  = Person is criminal

⇒ For population parameter, (Basic assumption) is

Null hypothesis ( $H_0$ )

⇒ The alternative hypothesis ( $H_1/H_a$ ) is opposite  
To Null hypothesis

→ Explanation :-

$H_0$  :- Null hypothesis

$H_1$  :- Alternate hypothesis / Research hypothesis.

Conditions

Present situation is best one.

•  $H_0$  :- Status quo /  $\rightarrow$  no difference / no action required

→ This correct :-  $\boxed{\text{No action}}$  required.

•  $H_1$  :- Researcher is Best / difference / action is b/w  $H_0 \& H_1$

→ This correct  
 $\boxed{\text{Take action}}$

Ex :-

$H_0$  : Person

is innocent (no action)

$H_1$  : Person

is criminal (action required)

### Example 2

Q :- A New Manufacturing method is believed to be better than the current method?

The new method is "No"

better than old method.

- A :-
- \* Null hypothesis [ $H_0$ ] :- better Than old method.  
(No action)
  - \* Alternative hypothesis [ $H_1$ ] :- **The new manufacturing method is better.**  
(Action required)

### Example 3

Feedback :-

\* SRK Sir  
(Data Science)  
is Better / Excellent?

\* Naresh I.T (Head)

Hypothesis Testing :-

15 students (response). \*\*\*\*\* = 5 Rating

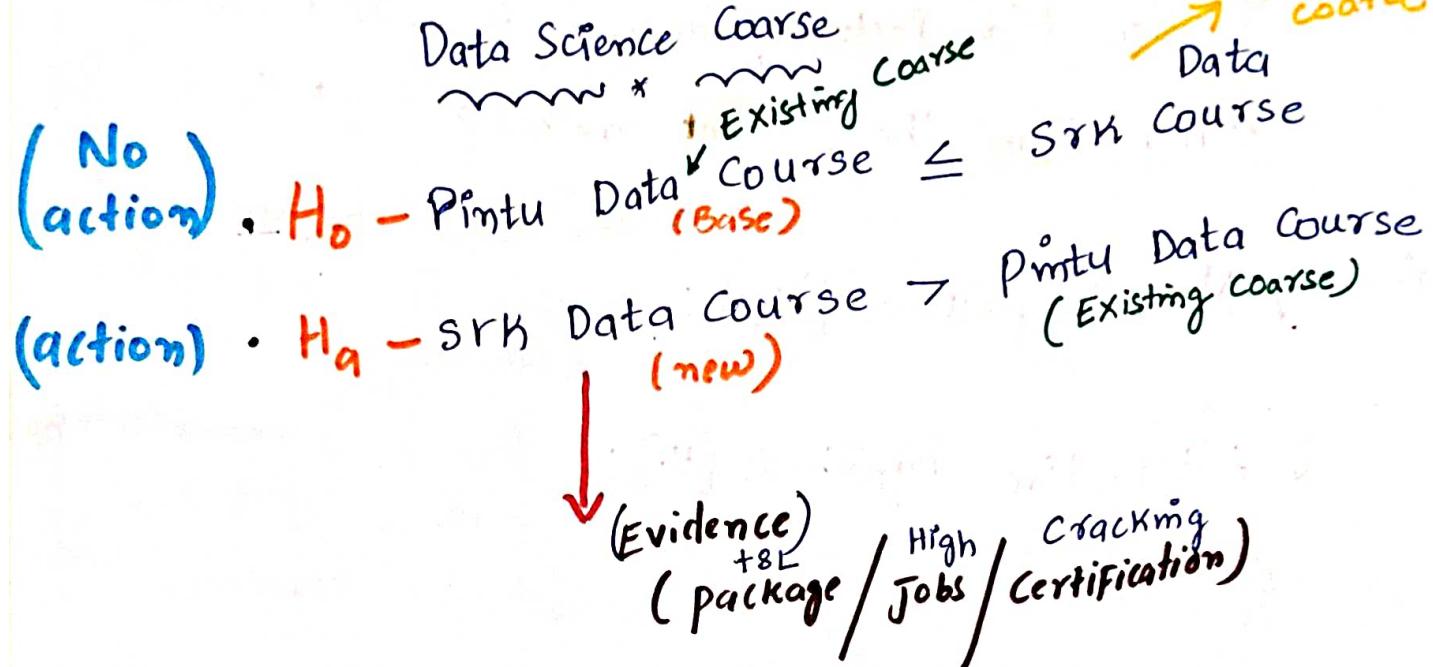
Less Than or Equal To 4.5/5 rating

$H_1 - \mu \leq 4.5$  (action required)

$H_0 - \mu > 4.5$  (No action required)

Example 4 :-

(No action)



Example 5 :-

Q :- A New bonus plan, That is developed in attempt to increase sales

A :- Alternative hypothesis [ $H_1$ ] :-

The new bonus plan increase Sales

\*  $H_A$  = (Sales new plan  $\rightarrow$  Sales Existing plan)

\*  $H_0$  = (Sales new plan  $\rightarrow$  Sales Existing plan)

⇒ Null hypothesis [ $H_0$ ] :-

If New Bonus plan doesn't increases Sales

Equal also  
Don't consider ( $H_0$ )

\*  $H_0$  = (Sales new plan  $\leq$  Sales Existing plan)

Example 6 :- It sample Data is misleading  
population Data = Error (wrong sample)

Q :- A new Drug is developed with the goal of lowering cholesterol - Level more than Existing drug.

A :-

\* Alternative hypothesis  $H_a$  :-

The new drug lowers cholesterol level more than Existing Drug.

\* Null hypothesis  $H_0$  :-

The new Drug doesn't lower cholesterol - level more than Existing Drug.

Example 7 :- The Label on Milk bottle states that it contains 1000 ml.



Null hypothesis ( $H_0$ ) :-

The Label is Correct  $H_0: \mu \geq 1000 \text{ ml}$

Alternative hypothesis ( $H_a$ )

The Label is incorrect  $H_a: \mu < 1000 \text{ ml}$

The Label is incorrect

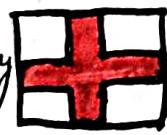
### \* Owner of milk

$$\cdot H_0: \bar{\mu} = 1000 \text{ ml} \quad (\text{No action})$$

$$\cdot H_1: \bar{\mu} \neq 1000 \text{ ml} \quad (\text{action})$$

### Example 8

Q:- A major hospital in chennai provides one of most comprehensive Emergency medical services in world.



- operating in a multiple hospital system with approximately 10 mobile medical units, The Service goal is to respond to medical Emergencies with mean time of 8 minutes (or) Less.
- The Director of medical Services wants to formulate a hypothesis test that could use a sample of times to determine whether Emergency response or not the Service goal of 8 minutes or

Less than is being achieved?

$$A:- \quad H_0: \bar{\mu} \leq 8 \text{ min} \quad (\text{No action})$$

$$(action) \Rightarrow H_1: \bar{\mu} > 8 \text{ min}$$

Mean of the vehicles

$H_0: \mu \leq 8$  [The Emergency Service is meeting the response goal; no follow up action required]

$H_1: \mu > 8$  [The Emergency Service is not meeting the response goal; appropriate follow-up action is necessary]

$\therefore \mu = \text{Mean response time of population of medical emergency requests.}$

⇒ Confusion Matrix

Example:- covid patient.

(Disease)

①

|                                |                         | Actual           |                         |
|--------------------------------|-------------------------|------------------|-------------------------|
|                                |                         | Person Has Covid | Person don't have Covid |
| Prediction (covid Test) Result | Person has Covid        | TRUE             | False (error)           |
|                                | Person don't have Covid | False (error)    | True                    |

②

|                             |                    | Person is Innocent | Person is Criminal |
|-----------------------------|--------------------|--------------------|--------------------|
| Prediction (Judge Evidence) | Person is Innocent | True               | False (Error)      |
|                             | Person is Criminal | False (Error)      | True               |

\* Confusion matrix

Ex : 3

$H_0$  = students understood class

$H_1$  = students didn't understand class

Actual (What happened)

|                                     |                           | Student understood | Student didn't understand |
|-------------------------------------|---------------------------|--------------------|---------------------------|
|                                     |                           | True               | False (Error)             |
| Prediction<br>(Feeling)<br>Evidence | Student understood        | True               | False (Error)             |
|                                     | Student didn't understand | False (Error)      | True                      |

For Example 8

Ardaan ayendki ani anukuntunna

Kaledhi ani anukuntunna

"Hospital Data"



[Actual]

|              |                          | Population Condition                   |                                        |
|--------------|--------------------------|----------------------------------------|----------------------------------------|
|              |                          | $H_0$ True ( $\mu \leq 8$ )            | $H_0$ False ( $\mu > 8$ )              |
| Conclusion   |                          | Correct Decision                       | Type I Error<br>(Incorrect Acceptance) |
| Accept $H_0$ | (Conclude $\mu \leq 8$ ) |                                        |                                        |
| Reject $H_0$ | (Conclude $\mu > 8$ )    | Type II Error<br>(Incorrect Rejection) | Correct Decision                       |

g/f/1  
30/3/22  
4:51 PM

Date:- 31/03/2022  
1:30pm

- \* The equity part of the hypothesis always appears in the null hypothesis.

Example :-  $H \geq H_0$

$H \leq H_0$

$H = H_0$

Mean

⇒ Example 9

Two Trainers.

Trainer 1

Already Existed

Trainer 2

New Trainer

$H_0$  :- Existing trainer

→ New Trainer

$H_1$  :- New Trainer

→ Already Existing Trainer

If both Trainer good ( $\Rightarrow$ )

We Take Null hypothesis

$H_0$  :- Existing trainer

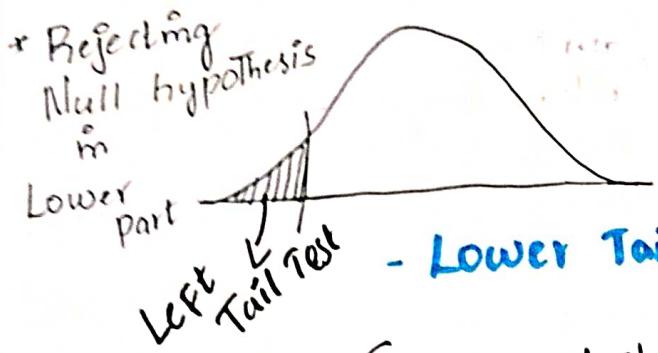
New Trainer

→  
Equals to

⇒ One Tailed [lower tailed] test

$$H_0 : \mu \geq \mu_0$$

$$H_1 : \mu < \mu_0$$

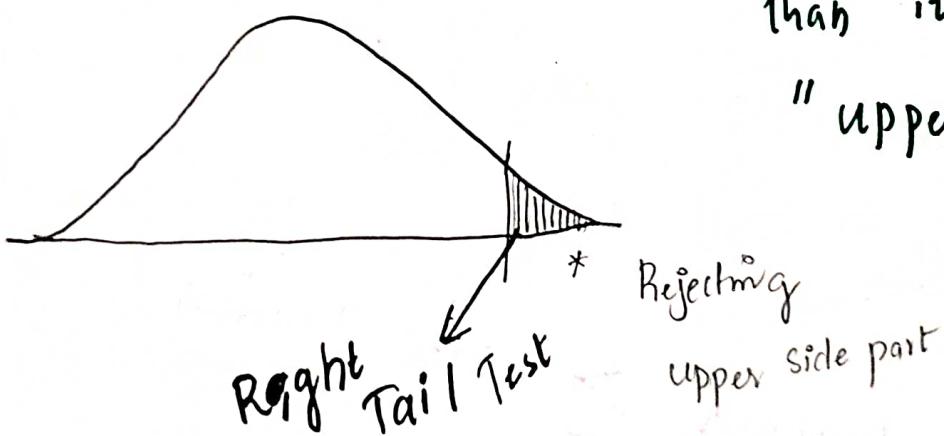


$H_1 / H_a$   
when ever The (Alternate) or Research Hypothesis Has  $<$  (Less than) Than it is called Symbol as "Lower Tailed Test."

⇒ One Tailed [upper tailed] test

$$H_0 : \mu \leq \mu_0$$

$$H_1 : \mu > \mu_0$$



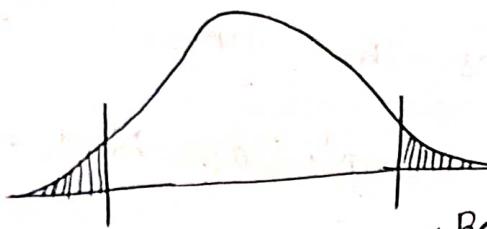
$H_1 / H_a$   
when ever the (Alternate) or Research Hypothesis Has  $>$  (Greater than) Symbol Than it is called as "upper tailed Test".

- Upper Tailed test.

\* Two Tailed test

$$H_0 : \mu = \mu_0$$

$$H_1 : \mu \neq \mu_0$$



$H_1 / H_a$   
when ever The (Alternate)  
Has ( $\neq$  not Equals to)  
Symbol, Then it is  
called as "two tailed  
Test".

### - Two Tailed Test

⇒ Steps of Hypothesis Testing

Step 1 :- Develop The Null and Alternative Hypothesis  
Identify

Step 2 :- Specify The Level of Significance " $\alpha$ "  
Here  $\alpha$  = Error Rate. To Accept.

$$\text{If } 5\% \text{ Error} = 0.05$$

$$10\% \text{ Error} = 0.1$$

$$20\% \text{ Error} = 0.2 \dots$$

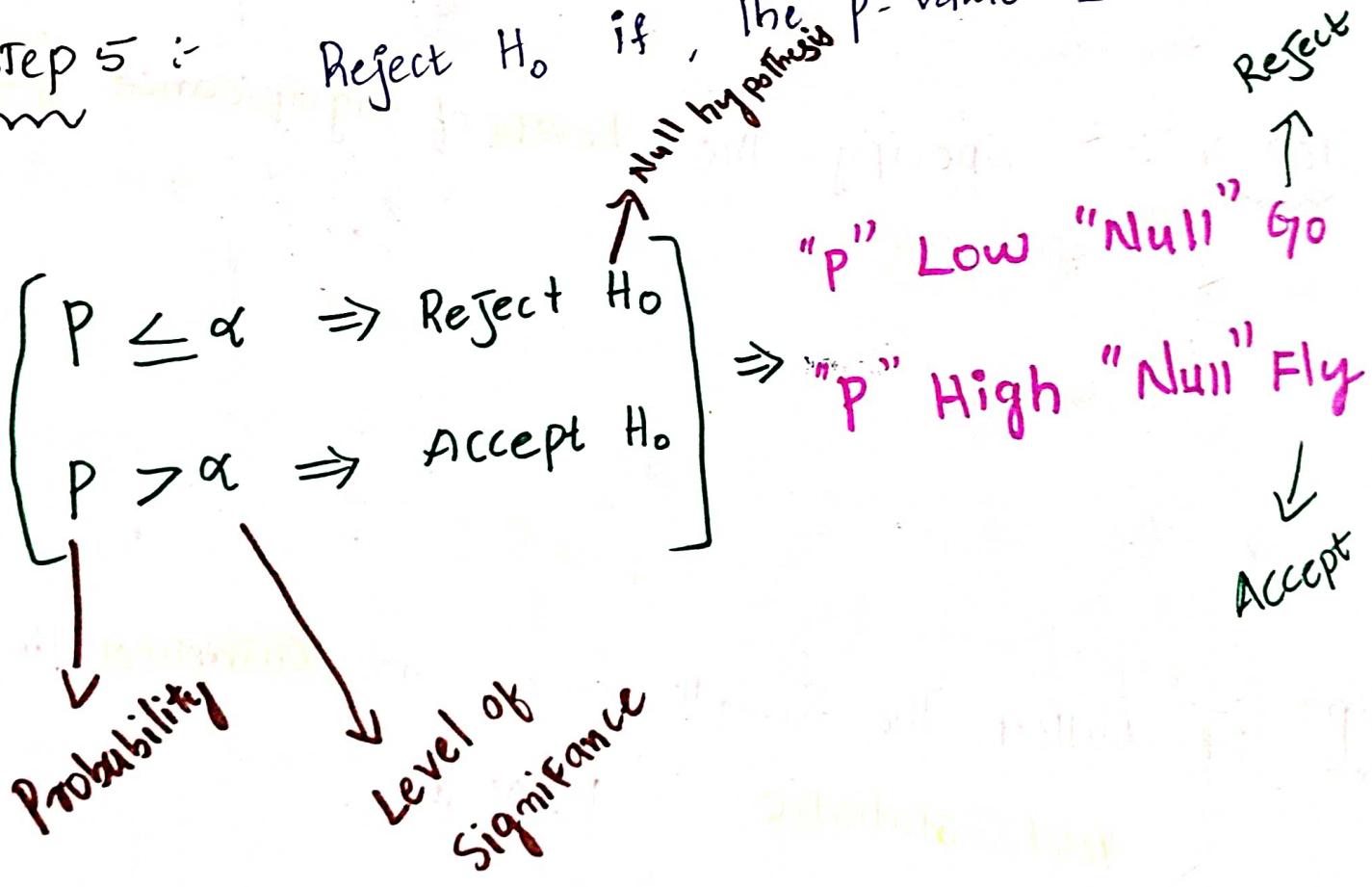
Step 3 :- Collect The Sample data and Compute The  
Test statistic "P-value".

Step 4 :- Compute The P-Value , Using Value of test static

\* The P-value is The Probability , Computed using Test statistic, that measures the Support (or Lack of Support) provided by sample of Null hypothesis.

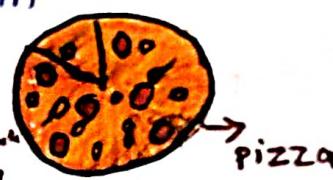
\* If the P-value is Less than (or) Equal to the Level of significance  $\alpha$ , the value of Test Static is in Rejection region.

Step 5 :- Reject  $H_0$  if , The P-value  $\leq \alpha$



\* Case Study  
One Tailed Tests about population when  $(\sigma)$  known.

Example 1

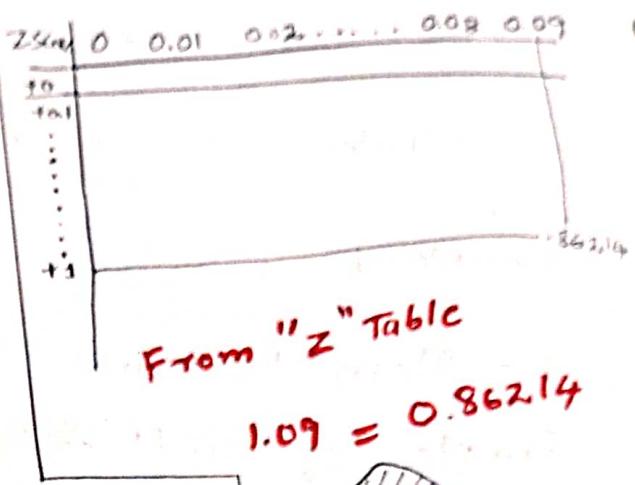
- The Mean Response Times For a Random Sample of "30 pizzas" Deliveries in 32 minutes. = Avg. val = 32, it is not meeting criteria  

- The Population Standard Deviation is believed to be 10 minutes
- The pizza delivery Services director wants to perform a hypothesis test, with  $\alpha = 0.05$  Level of Significance, to determine whether the Service goal of 30 minutes or less is being achieved.

- Answer :-  
 But, we are No. of Pizzas ( $n$ ) = 30  
 not working on Mean Sample ( $\bar{x}$ ) = 32 minutes  
 sample data. ( $\mu$ ) = 30 min  
 We have to estimate Standard deviation ( $\sigma$ ) = 10 min → Error  
 population. Level of Significance ( $\alpha$ ) = 0.05 Level  
 where, 30 minutes goal is So, 95% Accuracy  
 Achieved or not,  
 so, where, 30 = 32, we have to check with 95%.

Step 1 :-

$$H_0 : \mu \leq 30$$

$$H_1 : H > 30$$



Step 2 :-

Level of Significance = 0.05

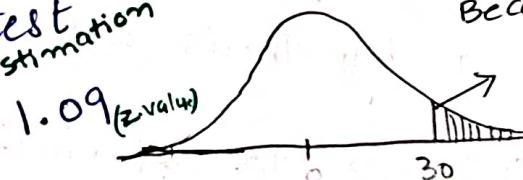
Step 3 :-

$$z_{\text{score}} = \frac{x - \mu}{\sigma}$$

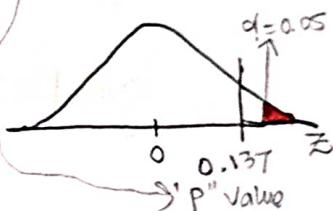
$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

10

$$\begin{aligned}
 & \text{Sample mean} \\
 & \text{Upper Tail test Estimation} \\
 & = \frac{32 - 30}{\sqrt{\frac{10}{30}}} = 1.09 \text{ (z-value)} \\
 & \text{Standard deviation} \\
 & \text{no. of samples} \\
 & \text{If } 1, 2, 3, 4 \\
 & \text{AS per}
 \end{aligned}$$



If Alternative hypothesis is  
 $H_1: M > 30$   
Because,   
greater  
called "upper tail test"



## " $\chi^2$ " - distribution

## "T" distribution

\* Population standard deviation  
 $(\sigma)$  is known

called "z" Distribution

If population standard deviation ( $\sigma$ ) is unknown  
Then it is called as  
"T" distribution.

$$Z = \frac{x - \mu}{\sigma} \rightarrow \text{Mean}$$

$Z$  distribution → Standard deviation

Ex :-

| Marks |
|-------|
| 90    |
| 60    |
| 58    |
| 83    |
| 98    |
| 32    |

(90)

$$\frac{x - \text{Avg}/\text{mean}}{\sigma} \rightarrow \text{Ex: } \frac{90 - 60}{10} = +3.$$

So, 90 is converted To  
Eqn.  $\rightarrow +3.$  by  $Z$  score  
from  $Z$ -Table

$$T = \frac{x - \mu}{\frac{s}{\sqrt{n-1}}} \rightarrow \text{Mean}$$

$n-1$  → degree of freedom

\* Bessel correction :

For Sample deviation, Sample Variance, we have to consider  $(n-1).$

⇒ Conversion of Every Value To "T" Score is called as

T Distribution

$$T_{\text{Test}} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

\* Step 4 :-

Probability of "T" Score

$$P(T \text{ score}) = \underline{\quad}$$

$Z$  distribution and T distribution are Equal

$$P(Z \text{ score}) =$$

\* Step 5

$$\text{If } n > 1000, Z = T$$

$$\text{Ex: } 0.738, 0.78 \rightarrow \text{Sample = 1200}$$

$Z \approx T$  up to 3,4 Decimals same answer

$$\text{If } n > 30, Z \approx T \approx 0.72, 0.73$$

up to First Digit Equal

Minimum Sample  
Should be  $30,$

Because,  $Z$  and  $T$

almost be "Equal"

## 'P" value approach

Step 4 :- Compute The P-value

For  $Z = 1.09$   
P-value = 0.137

Step 5 :- Determine whether To reject  $H_0$

Because,  $[P.\text{value} = 0.137 > \alpha = 0.05]$

Here

P.value  $>$  Level of significance  
 $(\alpha)$

So, We Accept "Null Hypothesis" ( $H_0$ )  
we Reject "Alternative Hypothesis" ( $H_1$ )

\* There are not sufficient statistical Evidence to  
infer that pizza Delivery Services is not  
meeting The response goal of "30 min".

## \* CODING

```
# import numpy as np
```

```
# from scipy import stats
```

- Null hypothesis = MeanTime  $\leq 30$

- Alternative hypothesis = Mean Time  $> 30$

- Alpha = 0.5

- Right Tail Z test  $\rightarrow$  Right Tail Rejection.

- Z calculated = 
$$\frac{\text{Observed mean} - \text{Population Mean}}{\text{Std. Error}}$$

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



```
# Z-test = (32 - 30) / (10 / np.sqrt(30))
```

```
print(Z-test)
```

```
[out] :-
```

1.0954451150

Normal Distribution

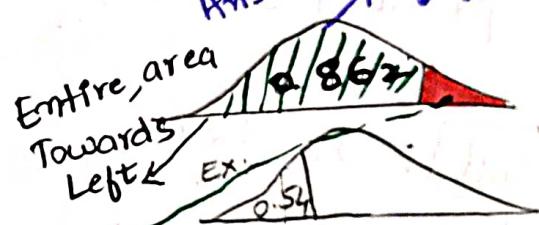
Cumulative Probability  
Distribution Function

Ans: Always  
 $\rightarrow$  left side

```
# stats.norm.cdf(1.09)
```

```
[out] :-
```

0.862143  $\rightarrow$  p value



For Right Tail Test

```
# 1 - stats.norm.cdf(1.09)
```

Out: 0.137  $\rightarrow$  Probability it is rejected

0.137  $>$  0.05

so, P High Null Fly/Accept)

Ex:-

If, Sample = 30, Mean Time = 28 minutes

so, it is satisfying the requirement  
Less Than 30. (No problem)

But,

If Sample = 30, Mean Time = 36

So, already we done with 32,  
We change 32 To 36. and calculate the

"z" Table, Z value,

$$\# z\_test = (36 - 30) / (10 / \text{mp.sqrt}(30))$$

print(z-test)

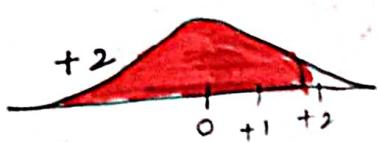
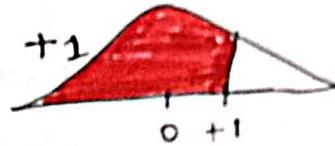
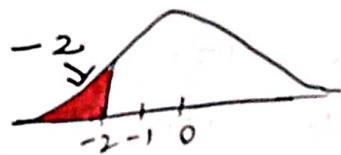
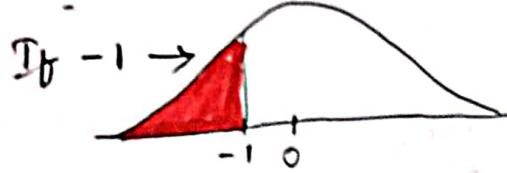
[out] :- 3.2826

$$\# 1 - \text{stats.norm.cdf}(z\_test)$$

[out] :- 0.00507

$\downarrow P < 0.05$   
 $0.005 < 0.05 \Rightarrow$  Reject Null Hypothesis.

\* Always towards left side



~~31/03/22 5.11pm~~

31/03/22  
9:30pm

### Case Study 2

#### Left Tail test (lower tail)

##### Example 2

A Machine is producing perfume bottles with the long term Average of  $150\text{ ml}$  and standard deviation of  $2\text{ ml}$ . Four bottles were picked and average volume was found to be  $149\text{ ml}$ . Has Volume  $\text{ml}$  of  $150\text{ ml}$  ? Test with

The average

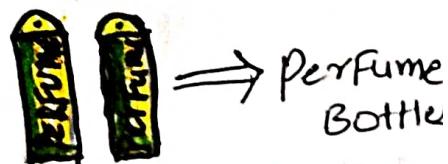
$95\%$  confidence level.

$$\text{Mean } (\mu) = 150 \text{ ml}$$

$$\text{Standard deviation } (\sigma) = 2 \text{ ml}$$

Ans:

$$n = 4$$



⇒ Perfume Bottles

Assume,  $n = 4$

$$\begin{array}{l} 1 \text{ bottle} = 151 \\ 2 \text{ bottle} = 150 \\ 3 \text{ bottle} = 148 \\ 4 \text{ bottle} = 149 \end{array}$$

$$\Rightarrow \begin{array}{l} \text{Avg} = 149 \\ \bar{X} = 149 \end{array}$$

\* Level of Confidence ( $\alpha$ ) = 0.05 (95%)

Step 1 :-

Alternate ( $H_1$ ) :-  $\mu < 150$  (action)

Null ( $H_0$ ) :  $\mu \geq 150$  (no action)

Step 2 :-



Step 3 :-

$$\alpha = 0.05$$

$$(1 - 0.95)$$

$$Z_{\text{Test}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{149 - 150}{\frac{2}{\sqrt{4}}} = \frac{-1}{1} = -1$$

standard deviation

no. of bottles

If less than 150, we have To reject

Step 4 :-

For Probability (-1)  $\Rightarrow$

From "Z" Table.

| $Z$  | 0       | 0.01 | 0.02 | 0.03 | 0.04 | 0.05 |
|------|---------|------|------|------|------|------|
| -0.1 |         |      |      |      |      |      |
| -0.9 |         |      |      |      |      |      |
| -1   | 0.15866 |      |      |      |      |      |

From  $Z$  score = 0.15866

$$P(-1) = 0.15866$$

Graphically



Step 5 :- So, it is Left Tail Test, we got Left Tail value as 0.15866, so we don't need To [1 - stats.norm.cdf(Z-test)]

$P(-1) = 0.15866 > 0.05$  "p" value is Greater Than " $\alpha$ " value

So, "p" High  $\rightarrow$  Null Fly (Accept The Null)

$\Rightarrow$  CODEING

Import numpy as np.

From scipy import stats.

#  $Z_{cal} = (149 - 150) / (2 / np.sqrt(4))$

$Z_{cal}$   $\rightarrow$  Z-Calculated Value.  
Out :- -1.0

# Left Tail  $\rightarrow$  stats.norm.cdf(z-cal)

# stats.norm.cdf(-1.0)

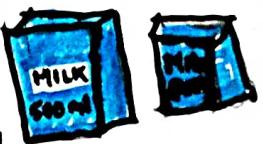
[out] :- 0.158655

### Case study 3

2 Tailed Test

Example 3 :- Milk Carton

- Assume that a Sample of 30 milk cartons provides a Sample mean of 505 ml
- Population standard Deviation 10 ml
- Perform Hypothesis test, at 0.03 Level of Significance, Population mean 500 ml. and to help determine whether the Filling process should Continue operating or be stopped and corrected.



Ans:-

$$* n = 30$$

$$* \text{Sample mean } (\bar{x}) = 505 \text{ ml}$$

$$* \text{standard deviation } (\sigma) = 10 \text{ ml}$$

- \* Level of significance ( $\alpha$ ) = 0.03
- \* population mean ( $H_0$ ) = 500ml

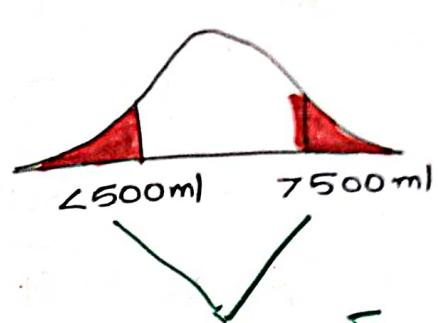
quality checked  
Not Less Than  
500ml

Step 1 :-

(owner perception)  
don't exceed  
more than 500 ml

$$H_0 : \mu = 500 \text{ mL}$$

$$H_1 : \mu \neq 500 \text{ mL}$$



Both The  
side Rejection

Step 2 :-

$$\alpha = 0.03 \text{ [Error]}$$

$$[1 - 0.97] = 97\% \text{ Accuracy.}$$

Step 3 :-

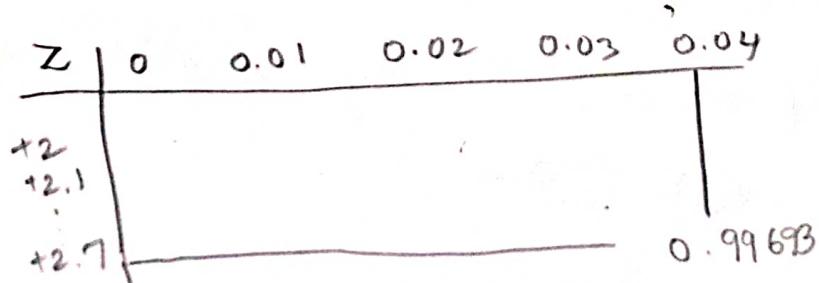
$$Z_{\text{Test}} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{505 - 500}{\frac{10}{\sqrt{30}}} = +2.7386$$

↑ sample mean  
↑ population mean / Estimate

Standard Deviation      no. of milk cartons

Step 4

$$P(2.7386) = \text{Let us consider it as } +2.74$$



$$P(Z < 2.7) = 0.99693 \rightarrow \text{Left Tail Test}$$

For "+" value Always do  
 $(1 - p\text{-value})$

So, for Right Tail Test

$$1 - 0.99693$$

$$= 0.003$$

If "z" value Have  
"-symbol"

Again, its 2 Tail Test

$$2 \times [1 - 0.9963]$$

$$P = 0.00616$$

$\frac{1}{2} \times (z\text{-score})$  Directly Multiply with "2"

calculator to vastha led  
But, Jupyter to  
vastkundi.

Step 5

$$P = 0.00616 < 0.03$$

so, "p" low, Null Go (Reject)

Rejecting Null hypothesis

→ There is no sufficient statistical Evidence to infer that

The Null hypothesis is True.

(i.e., The mean Filling quantity is not 500 ml)

## CODING

$$* \quad z_{\text{cal}} = (505 - 500) / (10 / \text{mp}.\sqrt{30})$$

print ( $z_{\text{cal}}$ )

[Out] :- 2.738612

$$\# \quad \text{Pvalue} = 2 * [1 - \text{stats.norm.cdf}(z_{\text{cal}})]$$

Pvalue

0.0061898

[Out]

DT: 1/04/22  
12:40 PM.

case study 4

enfj  
31/03/22  
11:30 pm.

## Fabric Data

### (1 Sample "Z" Test)

Question :- The Length of 25 Samples of a Fabric are Taken at random. Mean and Standard deviation From the historic 2 years study are 150 and 4 respectively. Test if the Current mean is greater than its historic mean, Assume  $\alpha$  to be 0.05 → If They don't give Data is Given in Excel sheet, Alpha value, Ans :- it has Only one column so, it is one sample.  $\alpha = 0.05$  (Standard Val)

$\alpha = 0.05$ ,  $\sigma = 4$ , so, it is "z" test.

Statistic Test

Parametric Test

(Entire given Data is normal)

Non parametric Test

\* Steps involved in Continuous Data (Hypothesis Testing)

Step 1 : Formulate,  $H_0$  and  $H_1$

Step 2 : Level of significance ( $\alpha$ )

Step 3 : check For Normality of given Data  
[For Continuous Data only]



For continuous data only we have Histogram not for Discrete Data.

Step 4 : Select Statistical Test and calculate "P" value

Step 5 : Based on "P" value, Conclude hypothesis Test.

\* How To Find "Normality" in The Data

• Option 1 : Skewness

• Option 2 : Density Curve

• Option 3 : Shapiro Test

### option 1 [skewness]

- If The skewness is between  $-1$  to  $+1$  [normal distribution]
  - If The skewness is Less Than  $-1$  [left skewed]
  - If The skewness is greater Than  $+1$  [right skewed]
- 

### option 2 [Density curve]

- If The density curve is Symmetrical [normal distribution]
- If the density curve is non Symmetrical [Skewed distribution]

\* option 3 [Shapiro test] For checking normality

$H_0$  : Data is normal

[if,  $P \geq 0.05$ , Normal]

$H_1$  : Data is not normal

[If  $P < 0.05$ , Not normal]

Ex:- Raymond Company.

Said and  
150cm (sold)  
Shirt Fabric Sample - 147 cm (But it having)

So, They want check samples.

\*  $H_0: \mu \geq 150$  [Fabric having 150 (or) more]  
No action

\*  $H_1: \mu < 150$  [Fabric having Less metres]  
action

Step 2

Level of significance  $\alpha = 0.05$

Step 3 check the normality.

# option 1 (skewness)

# Fabric[["Fabric\_Length"]].skew()

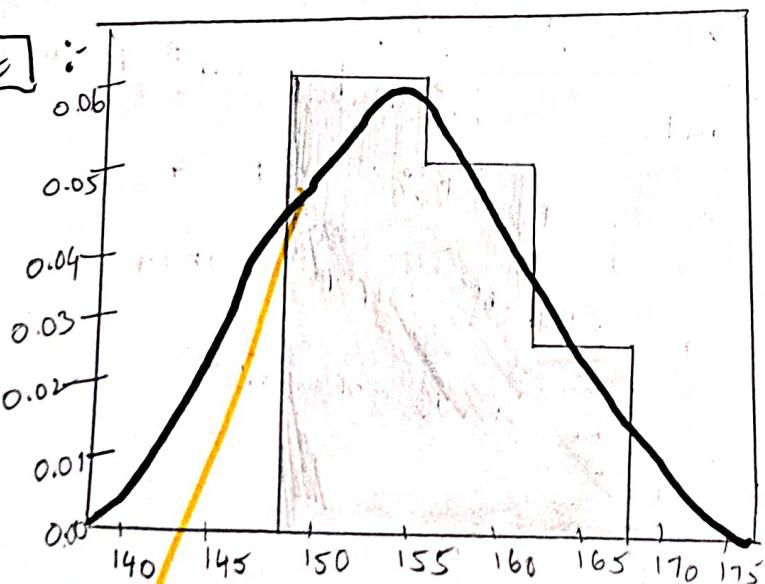
Out: -0.296

# option 2 (Density curve)

# sns.displot(fabric[["Fabric\_Length"]])

plt.show()

Out



So, we take this as "Last option" in find normal Distribution.

Here curve is Little Bended Fabric Length : How can we say it is "normal"

-1  
-0.1  
-0.2  
-0.3  
-0.4  
-0.5  
-0.6  
-0.7  
-0.8  
-0.9  
**1**

+1.0  
+0.1  
+0.2  
+0.3  
+0.4  
+0.5  
+0.6  
+0.7  
+0.8  
+0.9

# Option 3 [Shapiro Test]

# Stats.shapiro(fabric["Fabric\_Length"])

[Out] :- shapiro result (statistic = 0.9397..., pvalue = 0.1460)  
# p-value > 0.05, (Normal)

\* Doubt

Going Back To question

$$\mu = 150 \quad \text{Elakuda}$$
$$\mu \neq 150 \quad \text{Thesukovachis Kadha?}$$

But in question,  
we have Current Mean is "Greater than"

So, we take

$$H_0 - \mu \geq 150$$

$$H_1 - \mu < 150$$

Step 4

We have,

mean

To know  $\bar{x}$  value

$$Z_{\text{Test}} = \frac{\bar{x} - N}{\sigma / \sqrt{n}}$$

$$N = 150$$

$$\sigma = 4$$

$$n = 25$$

$$\bar{x} = ?$$

# fabric["Fabric\_Length"].mean()

[Out] :- 155.063

$$\# Z_{\text{Test}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$\# Z_{\text{Test}} = (155 - 150) / (4 / \text{mp.sqrt}(25))$$

Print  $Z_{\text{Test}}$

$$\boxed{\text{out}} := 6.325 \rightarrow Z_{\text{Test}} = P(6.325)$$

$$= P(0.999)$$

# stats.norm.cdf(z-test)

$$\boxed{\text{out}} := 0.999$$

Step 5 := Based on p.value, (Accept/Reject)

$$p = 0.999 > \alpha = 0.05$$

"P" High Null Fly (Accept)

So, we don't reject Null hypothesis.

Mean  $\geq 150$  (No action required)

To import :- Exceldata

fabric = pd.read\_excel("D://data science // shivaramakrishna// 5. inferential statistics // Fabric data.xlsx")

## Case Study 5

### BOLT Diameter (1 Sample T Test)

Example 5

The mean diameter of Bolt manufactured should be 10mm to be able to fit into nut. 20 Samples are taken at random from production line by a quality inspector. Conduct a test to check with 95% confidence that the mean is not different from the specific value.

Answer :-

Given, Data Set of 20 Samples.  
Here ( $\sigma$ ) is unknown so it is "T" test.

Step 1 :-

$$H_0 : \mu = 10 \text{ mm} \quad (\text{No action})$$

$$H_1 : \mu \neq 10 \text{ mm} \quad (\text{Action})$$

Step 2 :-

$$\alpha = 0.05$$

$$(1 - 0.95)$$

Step 3 :- Check the normality.  
# bolt [ "Diameter" ]. skew()

$$\boxed{\text{Out}} = 0.3809$$

Step 4 :- Statistical Test. ("1 sample" T Test)

$$T_{\text{test}} = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}}$$

mm  
वंदना  
veda?

In coding

# Stats. ttest\_1 Samp (bolt.Diameter, 10)

[Out] :- Ttest\_1 Samp Result (statistic = 4.114, p-value = 0.0005)

Step 5 :- Based on "p" value, Accept, Reject H<sub>0</sub>

$$P = 0.005 < \alpha = 0.05$$

so, "P" Low = Null go (Reject)

Reject "Null hypothesis" (avg. bolt diameter is not equal to 10) Take action.

1 sample = ~~Z test~~ T test (If n < 30, we have to Z test)

\* \* \* 1 sample = T Test (If n > 30, we take T test)

In Real life, we take T test mostly.

Example 6

## MARKETING STRATEGY

A Financial analyst at a Financial institute wants to evaluate a recent credit card promotion. After this promotion, 500 cardholders were randomly selected. Half received an ad promoting a Full Waiver of interest rate on purchases made over the next three months. The other half received a standard Christmas waiver. Did the ad promoting full interest rate increase purchases?

\* Understanding Question / with Different Example.

Ex:- FlipKart

365 Days

5% on axis  
Bank

(on Festivals)

10% Discount

on HDFC SBI

↳ So, it is giving good results

$$\times H_0 : 5\% \cdot (365) \geq 10\% \cdot \text{dis (Fstvl)}$$

$$\times H_1 : 5\% < 10\%$$

They are implementing every festival.

By doing Hypothesis Testing.

```

# import numpy as np
# import matplotlib.pyplot as plt
# import Seaborn as sns
# From Scipy import stats
# import pandas as pd

```

## Importing Data

```

# Promotion = pd.read_excel ("D:\\datascience\\Shiva rama
Krishna class\\5. Inferential statistics\\Promotion.xlsx")

```

Promotion.head()

2 sample Test.

Out :

| Interest Rate Waiver | Standard promotion |
|----------------------|--------------------|
| 1989.10              | 1272.35            |
| 1808.38              | 1250.38            |
| 1153.75              | 1474.78            |
| 1745.46              | 2064.89            |
| 1008.24              | 2030.87            |

Step 1:

$$\mu$$

$H_0$  : Avg. Purchases made by Full Interest Wave  $\leq$  Avg. Standard

(No action)

$H_1$  : Avg. purchase by FIW  $>$  Avg. Standard promotion

(action)

Step 2

Level of Significance =  $\alpha (0.05)$   
 $(1 - 0.95) \Rightarrow$

Step 3 (check Normality)  $\rightarrow$  Test For "2" Samples.

# print(stats.shapiro(Promotion. InterestRateWaiver))

[Out]: shapiroResult (statistic = 0.99), p value = 0.2245  
Greater Than 0.05  
So Normal Distribution

# print(stats.shapiro(Promotion. StandardPromotion))

[Out]: shapiroResult (statistic = 0.991, p value = 0.1915)  
Greater Than 0.05  
 $0.191 > 0.05$

### Variance Test / Levene Test

Fixed  $H_0$  :- Variances are Equal ✓  
 $H_1$  :- Variances are not Equal.

# stats.Levene(Promotion. InterestRateWaiver, Promotion.  
standardPromotion)

[Out]: LeveneResult (statistic = 1.13, p value = 0.287)  
Greater Than 0.05

Here, Variances are Equal,  
so, variances are Equal.

2 Samples are normal

We apply,

2 Sample Test For Equal Variance

Step 4 Select statistical Test & calculate 'P' value.

No ( $\sigma$ ) value, so, we take "T" Test

• 2 Sample "T" Test For Equal Variance.

# `stats.ttest_ind(promotion.InterestRateWaiver, promotion.StandardPromotion, equal_var=True)`

Out :- TTest Result (Statistics = 2.260, Pvalue = 0.024)  
P = 0.024 < 0.05

Step 5 :- Based on "p" value, Accept / Reject  $H_0$

$$P = 0.024 < 0.05$$

Here

"P" value is Less Than " $\alpha$ " value,  
"P" value is highly significant.

So,

"P" Low, Null Go (Reject)

"Full Interest"

$H_a$  = Avg. of purchases made by  
waiver is Greater than Avg.  
purchases made by standard  
promotion.

21/4/22 4:30pm

$Dt = g_{12} / g_{122} \left( \frac{12.51 \text{ nm}}{11+122 \text{ nm}} \right)$

Revising

Problem

$$H_0 : FIW \leq SC \quad \text{"2 sample T test"}$$

$$H_1 : FIW > SC \quad \rightarrow P < \alpha$$

(i) Levene Test

[Equal Variance]

$H_0$  : Variances are Equal

$H_1$  : Variances are not Equal

$\therefore P \geq \alpha \Rightarrow \text{Equal}$

$P \leq \alpha \Rightarrow \text{not equal}$

(ii) Shapiro Test

[Normality Test]

$H_0$  : Data is normal

$H_1$  : Data is not normal

$\therefore P > \alpha \Rightarrow \text{Normal}$

$P < \alpha \Rightarrow \text{not normal}$

Flowchart :-

2 Samples



Data is normal

↓ Yes

Equal Variance ?

↓ Yes

2 Sample "T" Test

For Equal Variance

No

Convert Skew To normal

No

2 Sample "T" Test for unequal Variance

## Case Study 7

[ANOVA/Ftest]

Example 7

### Contract Renewal

A marketing organization out sources their back-office operations to three different suppliers. The contracts are up for renewal and CMO wants to determine whether they should renew contracts with all suppliers (or) any specific supplier. CMO want to renew the contract of supplier with least transaction time. CMO will renew all contracts if the performance of all suppliers are similar.

From question →

1. All are Equal, if not
2. Less Time

Given Data,

3 Samples  
Supplier A      Supplier B      Supplier C

Ftest

← ANOVA Test

one way

Because, we have more than 2 samples in the data.

Import The Data.

Load Libraries:

```
# import numpy as np
# import pandas as pd
# import matplotlib.pyplot as plt
# import Seaborn as sns
# from Scipy import stats
```

Load the data

I always draw Data from where it stored.

So, I write Total address of the Data.

Mostly, we save in "PWD". But I store in another folder. So, I write Total address ✓

```
#Contract = pd.read_excel("D://data science // Shiwa Rama
                           Krishna // 5. inferential statistics // Contract Renewal
                           .xlsx")
```

contract.head()

out:-

| Supplier A | Supplier B | Supplier C |
|------------|------------|------------|
| 6.15       | 7.87       | 7.41       |
| 6.22       | 5.21       | 3.61       |
| 6.76       | 7.94       | 7.23       |
| 4.29       | 7.36       | 5.53       |
| 7.08       | 6.17       | 3.97       |

Step 1 :-

$H_0$  : Average Time by all Suppliers are Equal

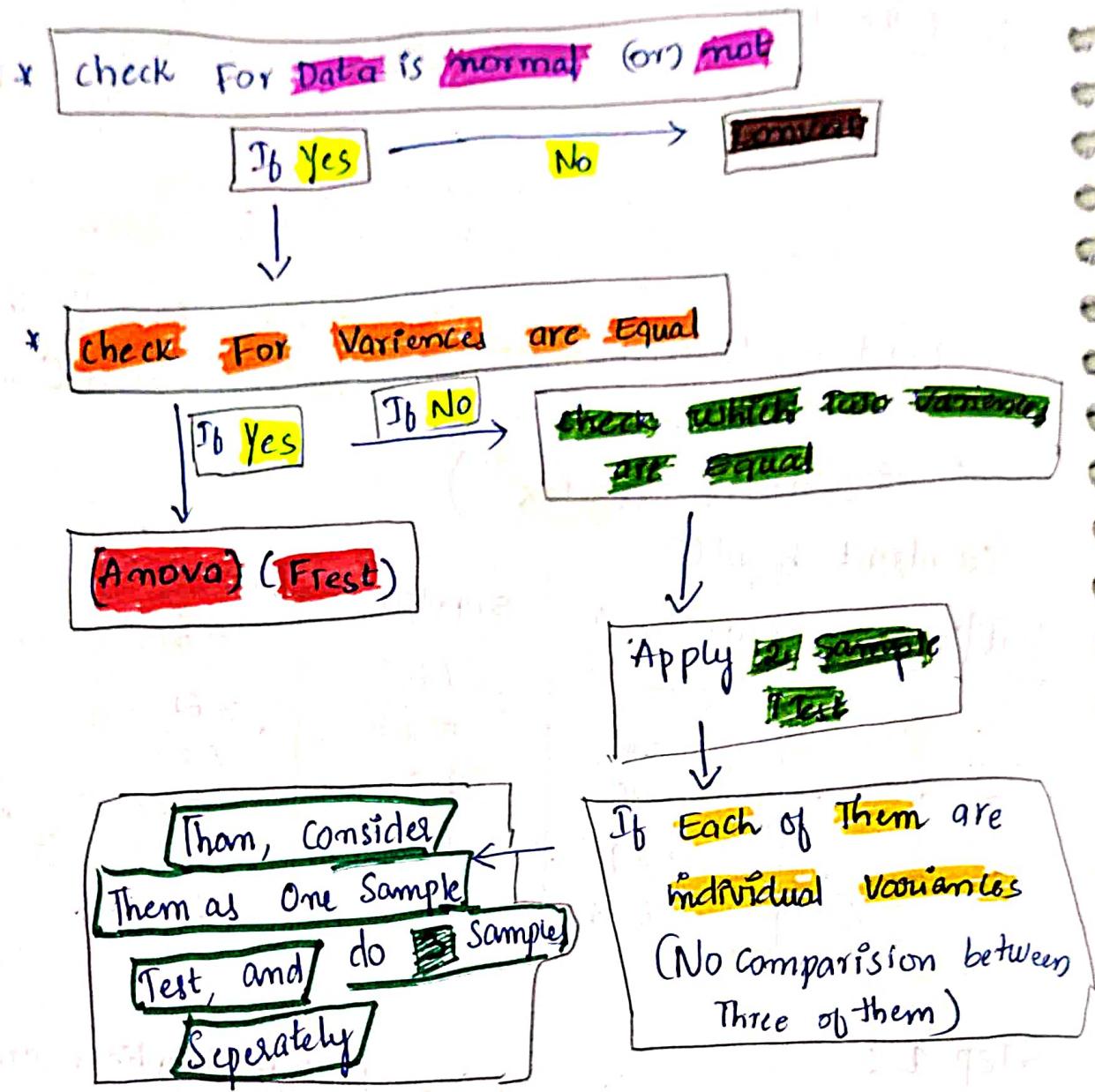
$H_1$  : Average Time by all Suppliers not Equal

Step 2 :-

Level of Significance  $\alpha = 0.05$

Step 3 :-

check for normality



Normality :-

```

# Print( stats.shapiro (contract. Supplier A))
, print( stats.shapiro . (contract. Supplier B))
, print( stats.shapiro . (Contract. Supplier C))
  
```

**Out**

```

Shapiro Result [ p.value = 0.89 ] normal
" [ p.value = 0.64 ]
" [ p.value = 0.57 ]
  
```

# Variance Test

# stats. Levene (contract. Supplier A, contract. Supplier B,  
contract. Supplier C) 3 inputs

out :- Levene Result ( pvalue = 0.77 )

P value > 0.05 (d)

Variences are Equal

Step 4 (Anova Test / F test)

# stats. f - OneWay (contract. Supplier A, contract.  
Supplier B, contract. Supplier C)

out :- F - OneWay Result ( pvalue = 0.103 )

P value > 0.25

From "p" value, p is greater than  $\alpha$ ,

'p' high, Null Fly (Accept)

We accept Null hypothesis.

= All the 3 suppliers have Equal mean

And, If we don't have 3 supplies <sup>Transaction Time</sup> Equal variance  
Then check, what a Sample Having Same variance

# print (stats. Levene (contract. Supplier A, contract. Supplier B))

# print (" " . " (contract. Supplier B, contract. Supplier C))

# print (" " . " (contract. Supplier C, contract. Supplier A))

## \* Discrete Variable

### Case study 1

#### Example 1 **FOOT ball Coach (1 proportion Test)**

The people carry out a poll to find the acceptability of a new football coach. It was decided that if the support rate for the coach for entire population was truly less than 25%, the coach would be fired - 2000 people participated and 482 people supported the new coach. Conduct a test to check if new coach should be fired with 95% level of confidence.

Sample people = 2000 ( $n$ )

question

Explanation:

$$n = 2000$$

(not Supporting) 1518

482 (supporting)

But, it is Sample  $\Rightarrow$  If less than 25% Fired.

May be in population

$$n = 10000 \leftarrow 7000$$

$$3000 = \frac{3000}{7000} = 30\%$$

$$\frac{482}{1518} = 24.1\%$$

From only we can conclude it is not reaching 25%.

## Population

|   |        |
|---|--------|
|   | Male   |
|   | Female |
| M |        |
| M |        |
| F |        |
| M |        |
| M |        |
| F |        |
| F |        |
| M |        |
| F |        |
| M |        |
| F |        |
| M |        |

$$m = 15$$

$$\text{Male} = 8$$

$$\text{Female} = 7$$

## Another Similar Example :

Sample Data

M, M, F, F

From the cam, Conclude

$$M = F (2 = 2)$$

By Sample Data,

Actual population Data is

$$M \neq F$$

$$8 \neq 7$$

Step 1 :-

$H_0$  : Coach not be Fired [No action] / No Difference

$H_1$  : coach to be Fired [action]

Step 2 :-

$$\alpha = 0.05$$

Step 3 :- check for normality / check Variance /  
we Don't have these in Discrete Data

Step 4 :- Select statistical Test, calculate "P" value.

\* 1 proportion Test

Because, we Have Only One column

# Stats. `bimom_test`

`Support` (4820, 2000, 0.25) → Total → Less than 25%  
`OUT`: 0.36615  
 "P" value 0.366 > 0.05  
 "P" High Null Fly (accept)  
 Donot reject null hypothesis.

(No data.xlsx given)

Directly written

Case Study 2  
 ~~~ \* ~~~

(2 proportion Test)

Example 2

Johnny Talkers

Johnny Talker Soft drinks division sales manager has been planning to launch a new sales incentive program for their sales executives. The sales executive felt that adults (>40 yrs) won't buy, children will & hence requested sales manager not to lunch analyze data & determine whether there is evidence at 5% significance level support the hypothesis

df = pd.read_excel ("JohnnyTalkers.xlsx")

df.head()

written
as it is
class
video
But,
Address
must
change.

[Out] :-

| | Person | Drinks |
|---|--------|------------------|
| 0 | Adults | did not purchase |
| 1 | Adults | did not purchase |
| 2 | Adults | did not purchase |
| 3 | Adults | did not purchase |
| 4 | Adults | did not purchase |

df ["Person"]. value_counts()

[Out]

children 740

Total = 1220

Adults 480

740 (L)

480 (A)

df ["Drinks"]. value_counts()

[Out] :-

didn't purchase 1010

Total
(1220)

1010 (d.p.)

purchase 210

210 (P)

CrossTable

pd.crosstab (df ["Person"], df ["Drinks"], margins=True)

[Out] :-

Drinks

Did not purchase

Purchased

All

Total

By, 1220.
are equal
(or) not
By seeing
Sample Data
we say
it is not
Equal,
But
what about

Person

→ Adults

→ Children

→ All

422

588

1010

58

152

210

480

740

1220

Here

Adults (purchased),

childrens (purchased)

$$\frac{58}{480} = 12\%$$

$$\frac{152}{740} = 20\%$$

population data? Equal (or) not

Step 1 :-

H_0 : Proportion of Adults \geq Proportion of Children

H_1 : proportion of adults $<$ proportion of children

Step 2

$$\alpha = 0.05$$

Step 3 :- No Normal / No variance — Discrete data

Step 4 :- 2 proportion Test

From statsmodels.stats.proportion import proportions_ztest

pur = np.array([58, 152]) # How many childrens are purchased

total = np.array([480, 740]) # Total no. of adults / children are there.

proportions_ztest(pur, total, alternative = "two-sided")

[Out] :- (-3.82, 0.00013)

$$p = 0.00013 < 0.05$$

Equal (or) not,
 $=, \neq,$

smaller

Larger
 $>$

Rufy
2/4/22

4:20 AM

Step 5 :- "p" Low Null Go (Reject)

| Person | Drinks | class |
|----------|-----------|--------------|
| Adult | Purchased | Low |
| Adult | not. P | middle |
| Adult | not. P | upper middle |
| children | not. P | High |
| adult | not. P | Low |
| children | Purchased | Low |

1 2 3 \Rightarrow 3 proportions

- In Person \rightarrow Adult
* Person \rightarrow children
- * Drinks \rightarrow purchased
* Drinks \rightarrow not purchased
- * class \rightarrow Low/middle/u.m
* class \rightarrow High.

Dt: 24/2/20
12:00pm

case study 3

(chi square Test)

Example 3

Bahaman Research

Baha Man Tech Research Company uses 4 regional centers in South Asia [India, China, Sri Lanka, Bangladesh] To input data of questionnaire responses. They audit a certain % of questionnaire response versus data entry. Any error in data entry renders it defective. The chief data scientist wants to check whether the defective % varies by country. help The manager draw appropriate inferences. ["1" means not defective, 0 means defective]

Load data set

~~* ~~

Bahaman = pd.read_excel("Bahaman.xlsx")

Bahaman.head()

Out:

| defective | Country |
|-----------|---------|
| 0 | India |
| 1 | India |

Count = pd.crosstab(Bahaman["Defective"], Bahaman["Count"])

Count

| Country
Defective | Bangladesh | China | India | Sri Lanka |
|----------------------|------------|-------|-------|-----------|
| 0 | 183 | 179 | 175 | 178 |
| 1 | 17 | 21 | 25 | 22 |

17/200, 21/200, 25/200, 22/200 = Equal (1) not equal (2)

Step 1:

Ho : No difference in proportion between Countries

H₁ : difference b/w countries

Step 2

Level of significance $\alpha = 0.05$

Step 3

No Normality / No Variance

Because it is
discrete data

Step 4

Chi-square Test (because we have more than 2 proportion
so, we take Chi-square Test in
Discrete data)

Contingency

stats. χ^2 - Contingency (count)

Cross Tab
Table.

out :- p value (0.63)

$p(0.63) > 0.05(\alpha)$

"p" high Null Fly (Accept)

We accept the hypothesis.

- Coding Completed ✓

Inferential
Statistics Completed

✓
X 2 4 2 2
3:06 PM