

Day - 1

Categorical column - integers of proportions

Continuous column - integers of mean & S.D.

Outlier -

$\text{abs } |x - \bar{x}| > 3\sigma$  is considered as

outlier - recommended.

treatment methods

i) Remove

ii) Replace with median

iii) Resampling

Extreme values - (outliers which are imp data)

treatment -

i) log transform (if skewed)

if very big range, then log transform  
compresses it.

ii) exponential - (if data is very narrow)

if the data is in a very narrow range  
then exponential transform spreads it.

## Descriptive stats -

gives some idea about the samples

collected. it doesn't talk about the

population.  $\rightarrow$  measures of central tendency

$\rightarrow$  measures of variance.

## Inferential stats -

we are doing one step more than descriptive

stats and are able to comment on

population as the sample is so-and-so

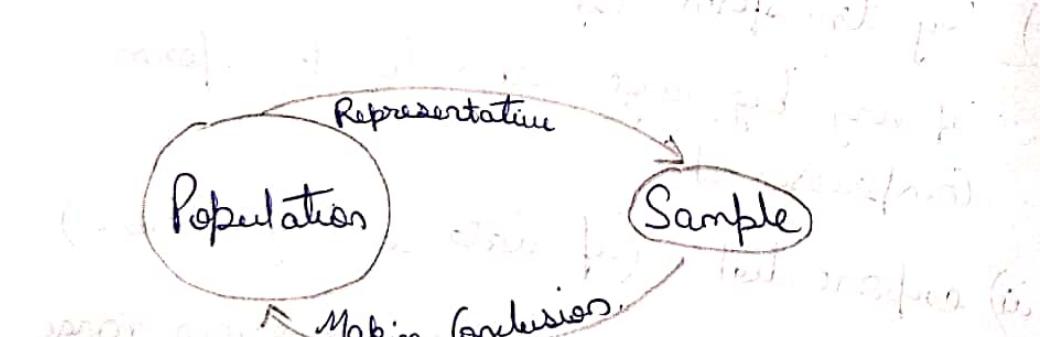
characteristics. Makes inference out of the

samples collected. From the samples

we are inferring the population.

t-test - for test of mean (continuous)

z-test - for test of properties (categorical)

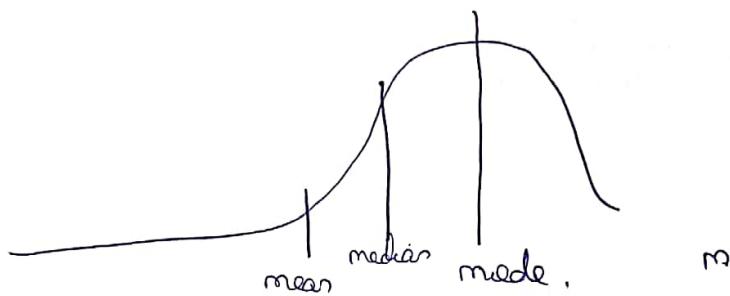


mode > median > mean.



Right.

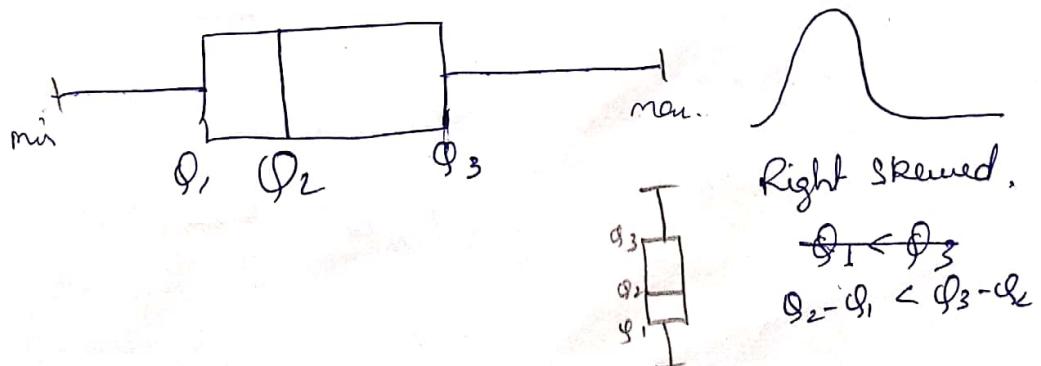
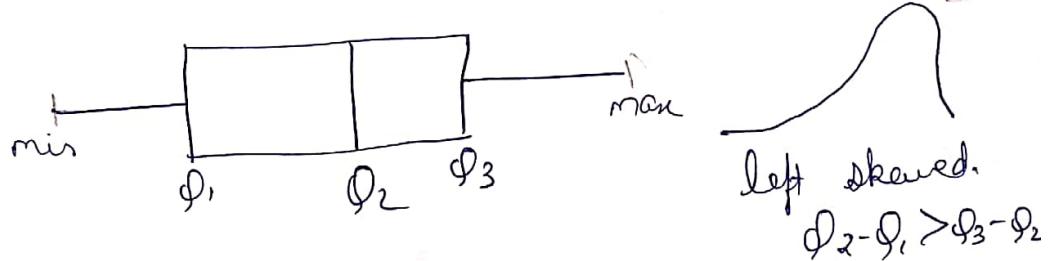
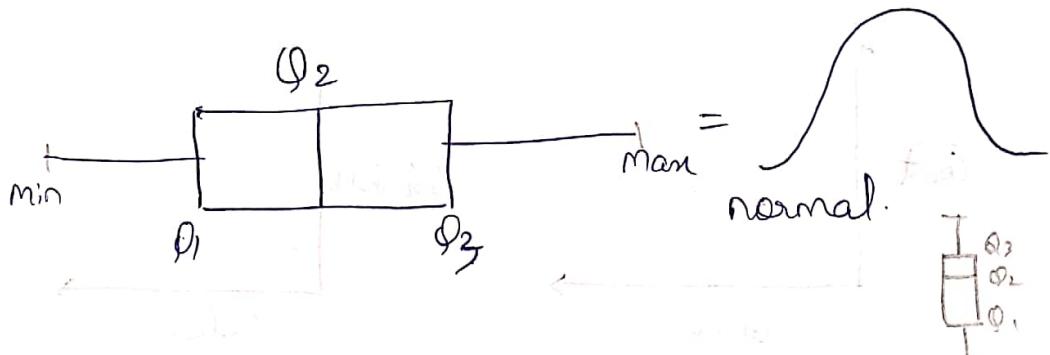
Skewed



Left

Skewed.

mean moves outwards because of extreme values.



## Linear model. (Parametric)

### Constraint

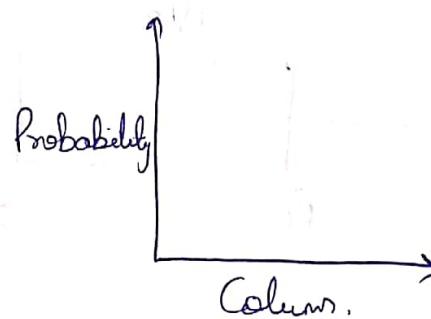
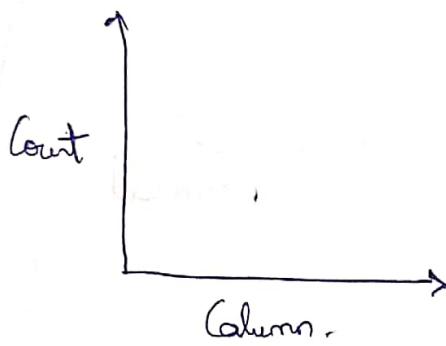
- Samples are random.
- Attributes follow normal distributions

## Non-linear models

Histogram

vs

Probability Density  
functions



\* Plot of all the probabilities is probability distribution

already known.

Probability -

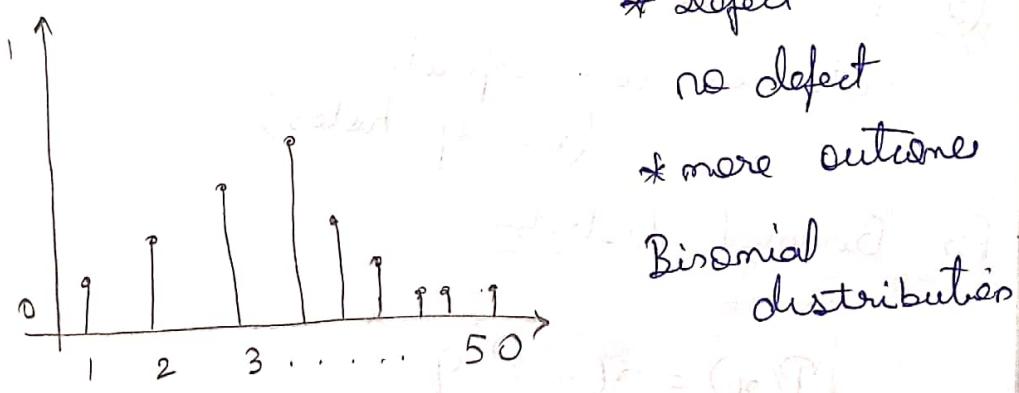
- 1) Classical probability  $\rightarrow$  Tossing a coin, dice, cards etc based on
- 2) Empirical prob.  $\rightarrow$  most used  $\rightarrow$  historical data.
- 3) Subjective prob. (not used)  $\rightarrow$  depends on person to person expert.

\* As the individual probability increases, the probability decreases.

## Probability distributions

① 5% defective

$n = 50$  t-shirt

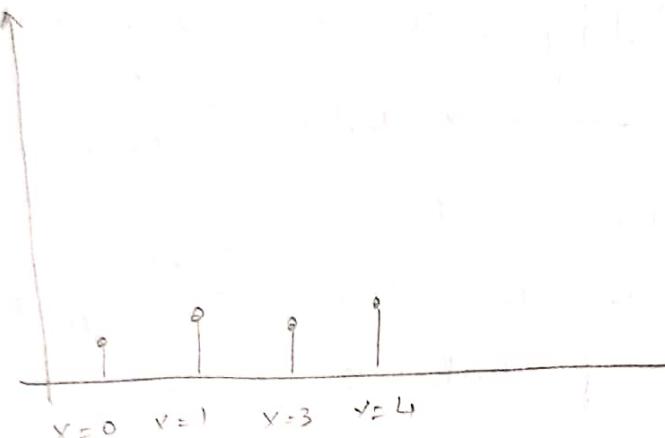


Cannot say 1.5 t-shirt different defect.  
only whole no's.

② Empirical data

→ avg defect is '7'

$n = 50$ .



← discrete →

many outcomes

Poisson

distribution

\* Difference b/w ① & ②.

① is Binomial distribution as limited no. of outcomes (defect, no defect)

② is Poisson's distribution as the outcomes are infinite (no. of holes)

for Binomial distrib-

$$P(x) = n(x) P^x q^{n-x}$$

$x = 1, n = 50, P = 0.05, \text{ so } q = 0.95$

$$P(x=1) = 50C_1 (0.05)(0.95)^{49}$$

For short time

for Poisson distributions -

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

e.g- T-shirt manufacturing data.

Process	Defects	Defects
1.	ND	0
2.	ND	0
3.	D	2
4.	ND	0
5.	D	3
6.		
7.		
:		
50khs.		

Setup  $\frac{150}{500000} = 0.03\%$

Two outcomes

more than two outcomes

for discrete values we use

Probability mass function (PMF)

Binomial  
Poisson.

and for Continuous values we use

Probability Density function (PDF)

normal

## Probability density functions -

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

\* characteristics of normal distributions -

→ Towards the mean value density is more, which implies, Average ~~is~~ always dominates.

A/C to central limit theorem - Average always dominates.

S.D for Binomial distributions -

$$S.D = \sqrt{np(1-p)}$$

$$np = 50 \times 0.15 = 7.5 \quad f \sim 8 \text{ with default.}$$

$$S.D = \sqrt{7.5(1-0.15)} \\ = \sqrt{6.375}$$

$$S.D. = \pm 2.52$$

Poisson distribution

$$S.D. = \sqrt{\lambda}$$

Normal distribution -

$$\text{mean.} = \frac{\sum_{i=1}^n x_i}{n}$$

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

$$\bar{x} = 24 \text{ years}$$

$$\bar{x} = 23.6$$

$$\sigma = 5 \text{ years}$$

$$\sigma = 2 \text{ years}$$

- spread is more

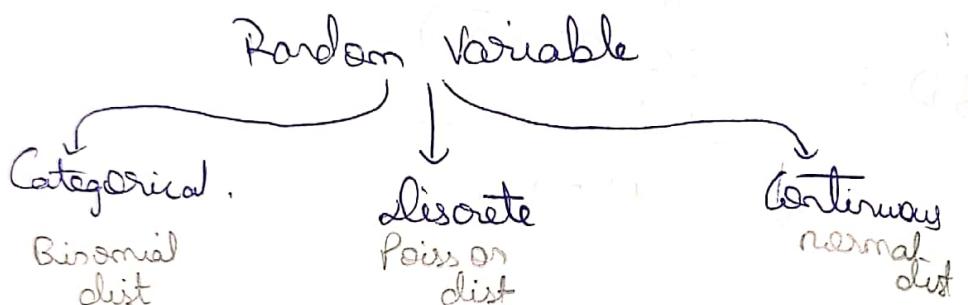
- spread is less

¶ On an average 6 customers every 1 min at a bank during 6z hours.

$$\lambda = 6.$$

what is the prob exactly 4 customers are given at a partic minute  
i.e  $x=4$  &  $\lambda=3$ .

### Day - 2



Categorical - This is also discrete in nature but it has generally only 2 categories - it is represented as Binomial e.g - 0/1, M/F, D/ND, Y/N etc, Tossing coin.

Discrete - This is multi category (more than 2)

e.g - no. of rooms in house. It is represented as Poisson distribution.

e.g - Throwing dice.

Continuous - This is of continuous values and follows normal distribution

e.g - Price of house etc.

→ In normal distributions if we calculate probability for one single point it will be 0.

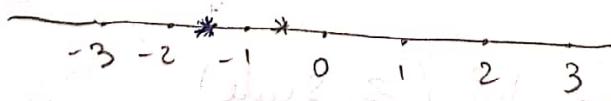
because as it is a continuous dist. we are calculating area under the curve. because at one single point we cannot integrate.

$$P(x = 45g) = 0. \text{ for normal dist.}$$

→ 
$$Z \text{ score} = \frac{x - \bar{x}}{S.D(x)}$$

all the values are mapped b/w -3 to 3.

→ In normal distributions, the values above the average  $\rightarrow$  right  
below the " "  $\rightarrow$  left  
as the 0 is mean in normal dist.



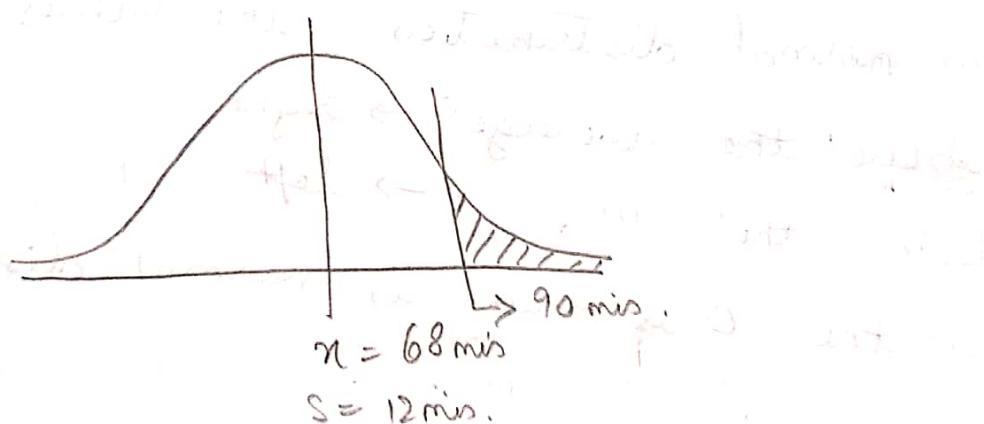
- \* Standard normal distribution -  
is a unitless probability distribution  
of the continuous variable.

why?  $\rightarrow$  we can use same scale to analyse  
different continuous variables.

- \* The norm. fun. always calculates area  
on the L.H. S.

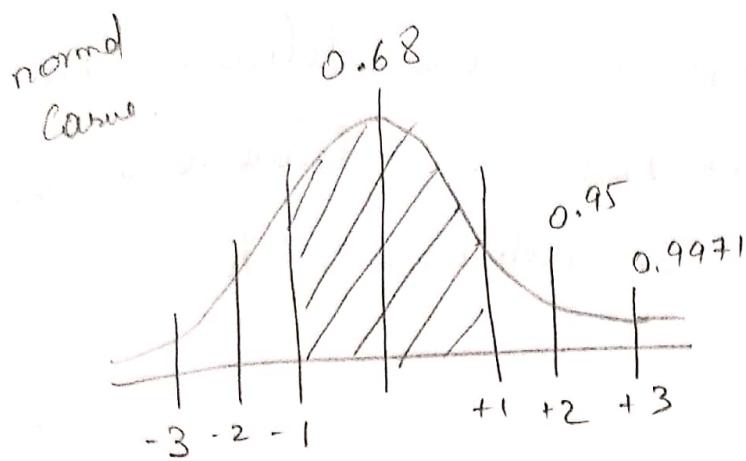
Q - a Survey Conducted, found average 68 min  
on phone, S.d is 12 min.

a) what proportion of S.P users are spend  
more than 90 min?



first convert 90 to z scale.

$$= \frac{90 - 68}{12} = 1.8 \text{ (z scale)}$$

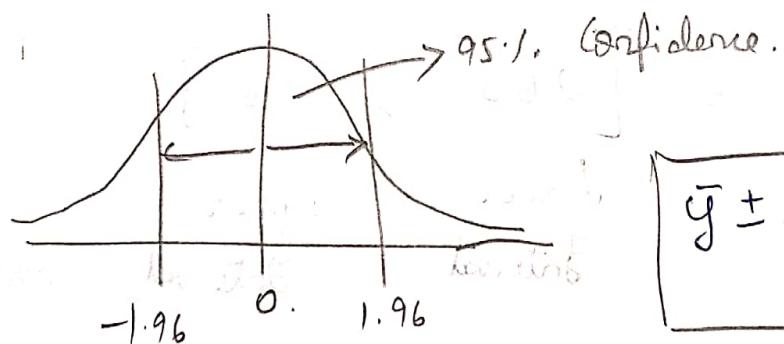


\* Information is contained in a spread.

- 30 to 30.

as the spread of the interval increases our confidence that the inference is correct for the sample to population from sample will also increase.

$$\text{Standard error} = \frac{s}{\sqrt{n}}$$



$$\bar{y} \pm \frac{1.96 \times s}{\sqrt{n}}$$

Q - Courier Company, mean delivery less than 3 years. 50 deliveries random, determines mean delivery = 2.8 hrs,

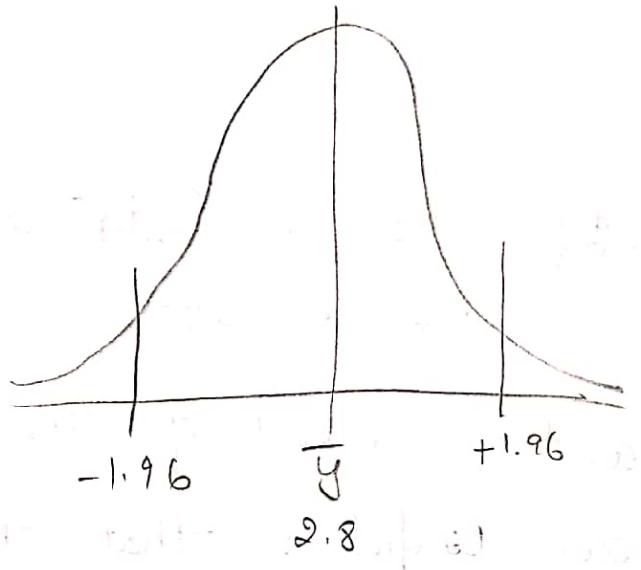
$$s = .6 \text{ hrs.}$$

$$\alpha = 95\%$$

$$n = 50$$

$$\bar{y} = 2.8 \text{ hrs.}$$

$$s = 0.6 \text{ hrs.}$$



$$C.I = \bar{y} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$

$$= \left[ 2.8 - 1.96 \times \frac{0.6}{\sqrt{50}}, \quad 2.8 + 1.96 \times \frac{0.6}{\sqrt{50}} \right]$$

$$= [2.63, 2.96]$$

Lower Interval      Upper Interval

→ The interval from 2.63 to 2.96 form a 95% Confidence interval for mean( $\mu$ ). In other words we are 95% confident that the average delivery time lies between 2.63 & 2.96 hrs.

Confidence Coefficient	Value of $\alpha/2$	Area in Table 1 $1-\alpha/2$	Corresponding z-value $z_{\alpha/2}$
0.90	0.05	0.95	1.645
0.95	0.025	0.975	1.96
0.98	0.01	0.99	2.33
0.99	0.005	0.995	2.58

$$\begin{array}{cccc}
 90\% & 95\% & 98\% & 99\% \\
 1.645 & 1.96 & 2.33 & 2.58
 \end{array}$$

→ when Sample Size increase, it becomes a very good representation of the populations.

← Instead of  $\alpha = 5\%$  use  $\alpha = 1\%$ .

$$= \left[ 2.8 - \frac{2.58 \times (0.6)}{\sqrt{50}}, 2.8 + \frac{2.58 \times (0.6)}{\sqrt{50}} \right]$$

$$= [2.54, 3.01]$$

## Session - 3

$$1 \text{ min} = 0.40 \\ 120 \text{ min} = ?$$

$$n = 80 \text{ min.}$$

1 hr 20 min.

$$n = 100$$

$$\bar{y} / \bar{x} = 4.5$$

$$S = 1.2$$

$$C.I = \bar{y} \pm 1.96 \times \frac{\sigma}{\sqrt{n}}$$

$$= \left[ 4.5 + 1.96 \times \frac{1.2}{\sqrt{100}}, \quad 4.5 - \frac{1.96 \times 1.2}{\sqrt{100}} \right]$$

$$= [4.735, 4.2648]$$

If  $n > 30$

Then Sample S.D  $\approx$  Population S.D.

To Calculate C.I -

LCI, UCI = stats.norm.interval( $(C_i, \text{loc} = \bar{x}_{\text{avg}})$ ,  
Scale =  $s$ )

$$C_i = 0.95$$

$\bar{x}_{\text{avg}}$  = mean.

$$S = \frac{\sigma}{\sqrt{n}} \quad \text{S.E Standard error.}$$

## t-distributions -

→ adjusted Std. dist., used for hypotheses testing

$$\frac{x - \bar{x}}{s/\sqrt{n}}$$

P- Hindustan Pencils; Pencil mean length

172 mm, S.d 0.02 mm.

Sample taken,  $n = 100$ ,  $\bar{x} = 170$  mm

$\alpha = 95\%$ .

$x = 172$  mm.

$$C.I = \bar{x} \pm 1.96 \times \frac{s}{\sqrt{n}}$$

$$= 170 \pm 1.96 \times \frac{0.02}{\sqrt{100}}$$

$$= \left[ 170 - 1.96 \times \frac{0.02}{\sqrt{100}}, 170 + 1.96 \times \frac{0.02}{\sqrt{100}} \right]$$

$$= [169.99, 170.003]$$

as the ~~std. deviation~~ mean 172 is not in the sample confidence interval, the sample is not inferential for populations this means, The S.D is little higher than 0.02,

## Hypothesis testing

→ If the claim is alternate hypothesis is on continuous variable, it will be on near. (Categorical)

If the claim is on discrete variable, it will be on the properties.

Continuous → near.

Categorical → Properties.

$H_a: \mu_0 < 3 \text{ hrs}$  (Couriers Company of India Ltd. (G))

$H_0: \mu_0 \geq 3 \text{ hrs}$ .

as a  
data Sci  
always  
contradict

soft  
inequality

One  $< >$  - Hard inequality

Tailed test.  $\leq \geq$  - Soft inequality.

(for null hypothesis).

Modelling Hypothesis  $\rightarrow$  for modelling we only deal with = or  $\neq$

$H_0: =$   
 $H_a: \neq$

= or  $\neq$  two tailed test  
 $>< \rightarrow$  one tailed test

In inferential the soft inequality must be given to null hypothesis

whereas in modelling hypothesis, the equality should be given to null hyp.

Inferential,  $\# H_0: \mu_0 \geq 3\text{hr}$   
 (one tailed test)  
 $H_a: \mu_0 < 3\text{hr.}$

Modelling  $H_0: \mu_1 = \mu_2$  (Two tailed test)  
 $H_a: \mu_1 \neq \mu_2$

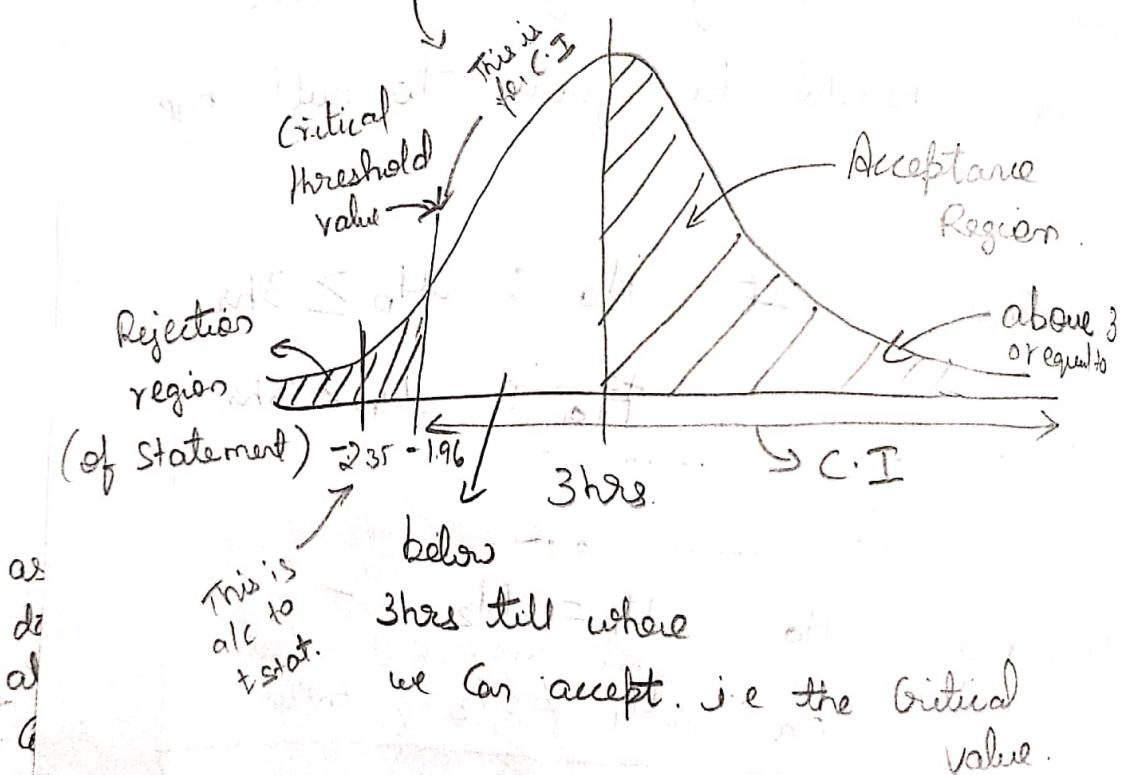
\* always start with alternate hypothesis

→ Carrier problem, where they had stated that the <sup>mean</sup> delivery time is less than 3 hrs.

$$H_a : \mu_0 < 3 \text{ hrs.}$$

$$H_0 : \mu_0 \geq 3 \text{ hrs.}$$

graphically.

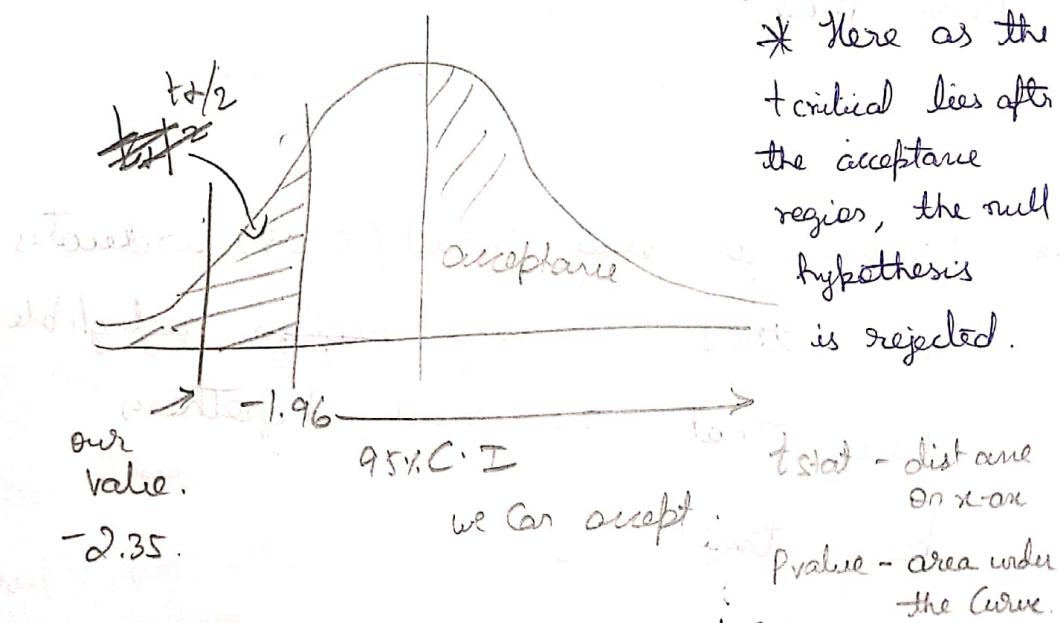


\* One Sample t test,  $\Rightarrow$  left tail test.

$$\begin{aligned}
 \text{(critical value)} \Rightarrow t_{\text{stat}} &= \frac{\bar{y} - 3}{s/\sqrt{n}} \\
 &= \frac{2.8 - 3}{0.6/\sqrt{50}} = -2.35
 \end{aligned}$$

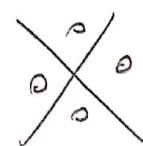
7 The left H.S. side is contradicting for the null hypothesis, but not the whole of L.H.S is contradicting, we are ready to accept till one point. i.e 95%. we are accepting. Only 5%. we are reserving for the rejections.

\* For 95% C.I., minimum  $-1.96$  the acceptance should lie away.



If  $t_{\text{stat}}$  lies within the  $-1.96$  critical value is  $t_{\text{stat}}$ . and the area w.r.t to critical value it is the p-value.

$t_{\text{stat}} \rightarrow$  distance on x-axis  
 $p\text{-value} \rightarrow$  area under curve.



## Rule for Rejecting null hypothesis.

95%, 99%

$t_{\text{stat}} > 1.96$  or  $t_{\text{stat}} > 2.58$

or  $P\text{value} < 0.05$  or  $P\text{value} < 0.01$

$p$ -value is a probability of null hypothesis being true

\*  $P\text{value}$  is very small (0.05) indicates that there is very negligible chance that the null hypothesis being true.

\* 5% error, the rule or  $\alpha$  is relaxed as compared to 1% error which is very strict.

at what level of confidence  $\rightarrow$  client will provide.

Continued - delivery Company example -

Calculate p-value for t-stat  $S_{\bar{x}} = 3.35$

$$\text{Stats.norm. Cdf}(-2.35) = 0.0092 \rightarrow \text{p-value.}$$

which is less than 0.05. This function uses normal dist. for accurate p-value, use t dist.

$\therefore$  There is less than 5% chance for the null hypothesis being true. Here the remaining 95% supports the claim made by the delivery company.

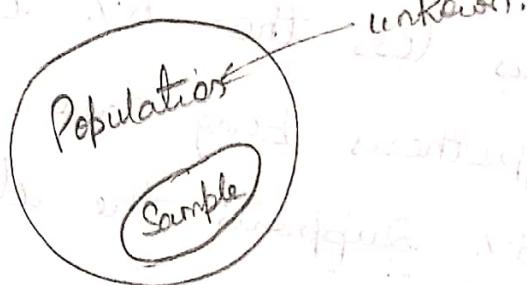
$\therefore$  We fail to accept the null hypothesis.

So acc to the alternate hypothesis, the claim made by the delivery company is true i.e. they make all their deliveries within 3 hours.

\* To get accurate p-value, we should use the t-distribution. In the previous case we are using the normal distribution.

t test → test of mean

One Sample →



Two Sample t test →

\* Comparing two groups

$$\begin{array}{c} n_1 = 1000 \\ \bar{y}_1 \\ s_1 \end{array} \quad \begin{array}{c} n_2 = 2000 \\ \bar{y}_2 \\ s_2 \end{array}$$

$$t_{\text{stat}} = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

when

One Sample t-test has 3 Case -

- 1)  $\leq$  left tail
- 2) Right tail  $>$
- 3) Two tail  $=$

Two Sample t-test has only one -

- 1) Equal to  $=$
- 2) Not equal to  $\neq$

In two Sample t-test -

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

→ Functions for t-test →

{ for t<sub>0.05</sub> critical value -

→ t-test - 1 Samp.

→ we need Samples So we have to generate Samples with random values

\* randn → will gen (0 to 1),  $\mu = 0, \text{SD} = 1$

\* randint → integers

\* If we want to generate a Sample with required S.D and mean  
we just have to reverse the Z val

$$Z_{\text{val}} = \frac{x - \bar{x}}{\text{S.D}(x)}$$

$$Z_{\text{val}} \times \text{S.D}(x) + \bar{x} = x.$$

S.D of x.

for ~~delivery~~ HR example  $\rightarrow$

$$\bar{x} = 2.8 \quad \mu = 3$$

$$S.D = 0.6$$

$$n = 50$$

for generated Sample is jupyter.

$$\mu = 3$$

$$\bar{x} = 2.81$$

$$S = 0.61$$

Python further to Cal. t statistic Sample t-test

$$t\text{test\_stat, pval} = t\text{test\_1samp}(S\_std, 3)$$

\* This will give us t statistics and p-value.

Sample population mean.

Q Sayabean e.g -

$$\bar{y} = 520 \quad 573 \quad \mu = 520$$

$$S = 124$$

$$n = 70$$

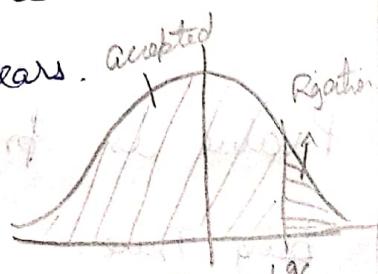
Company is claiming yield is greater than

520. Compared to last 2 years.

$$H_a : \mu > 250$$

$$H_0 : \mu \leq 250$$

$$t\text{stat} = 3.41 \quad P\text{val} = 0.001$$



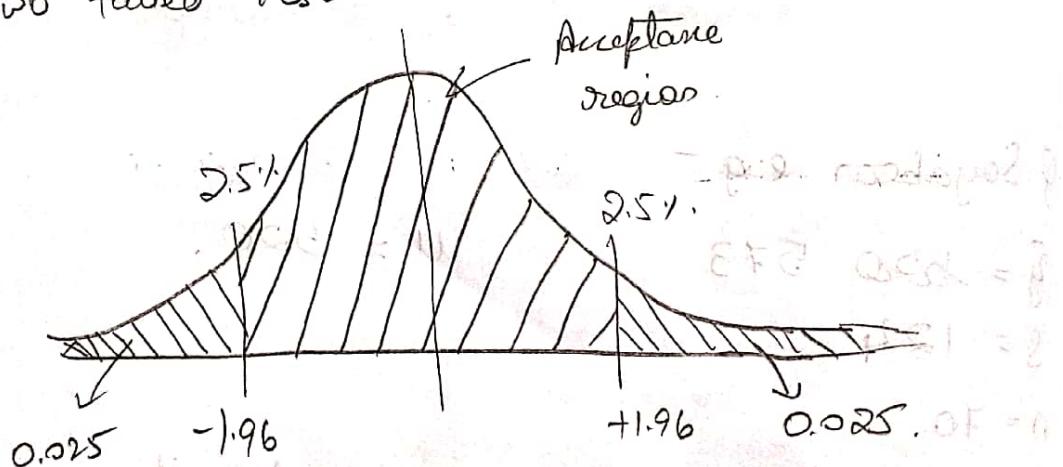
∴ as  $t_{\text{stat}} > 1.96$  and  $P_{\text{val}} < 0.05$   
 we reject null hypothesis.

error - To reject  $H_0 \rightarrow$

Pval.	$\alpha$ Pvalue <	$\alpha/2$	tstat.
90%	0.1	0.05	1.64
95%	0.05	0.025	1.96
98%	0.02	0.01	2.33
99%	0.01	0.005	2.58

Remember

Two tailed test -



Pvalue is probability of null hypothesis being True.

\* If mean and median values are lying closer, the data is normal.

class data -

$$n = 30$$

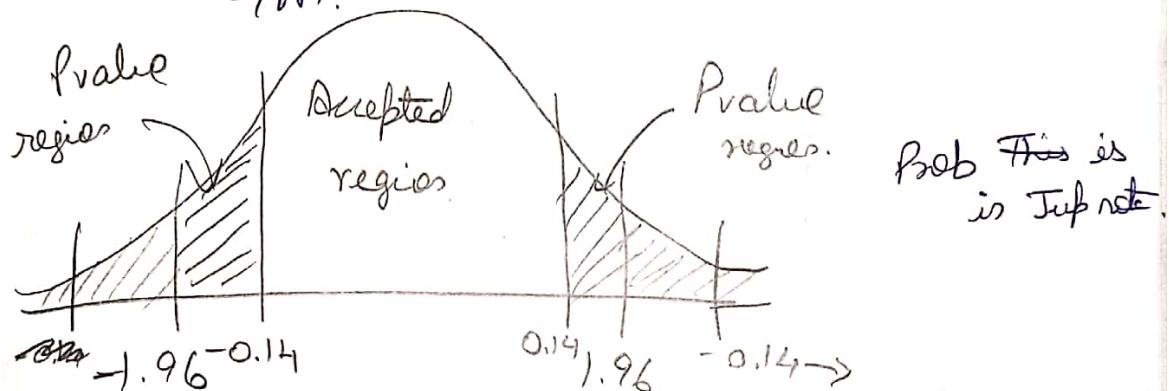
$$\bar{x} = 24.95$$

$$s = 2.48$$

$$H_0: \mu = 25 \text{ yrs}$$

$$H_a: \mu \neq 25 \text{ yrs.}$$

$$t_{\text{stat}} = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{24.95 - 25}{2.48/\sqrt{30}} = -0.1424$$



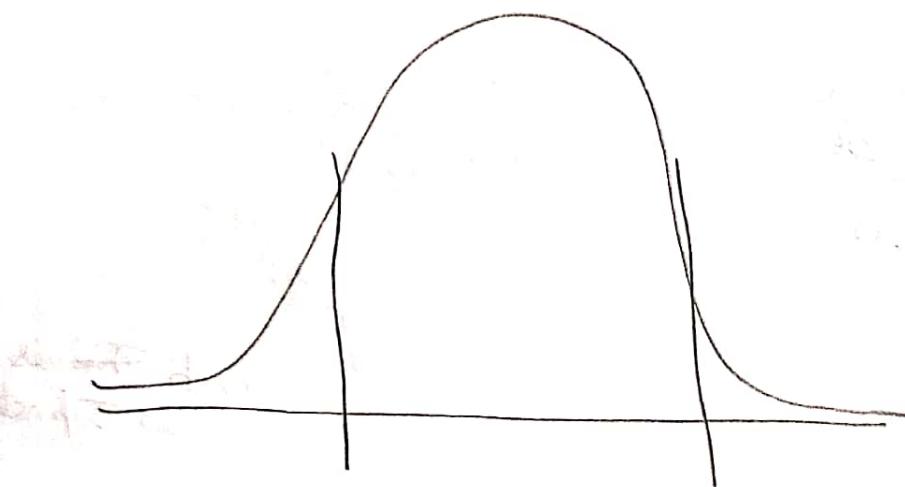
Here  $0.14 < 1.96$  and Pvalue is 0.88.

This means that There is 88% chance of  $H_0$  being true.

So we accept  $H_0$ .

Therefore, the class age is average of 25 yrs.

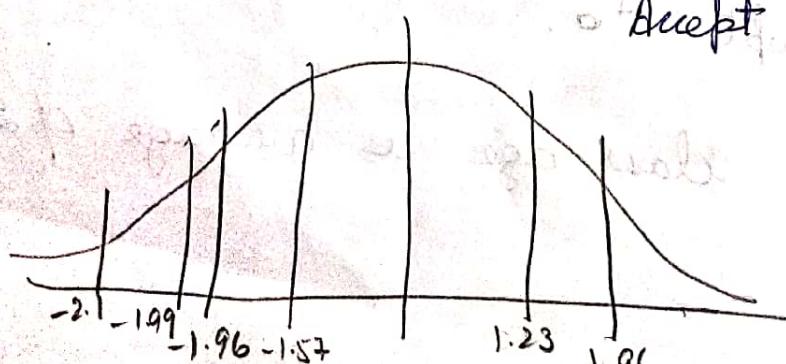
→ In the two Sided t-test, we have to check the error on both sides as we are asking whether it is equal to or not equal to.



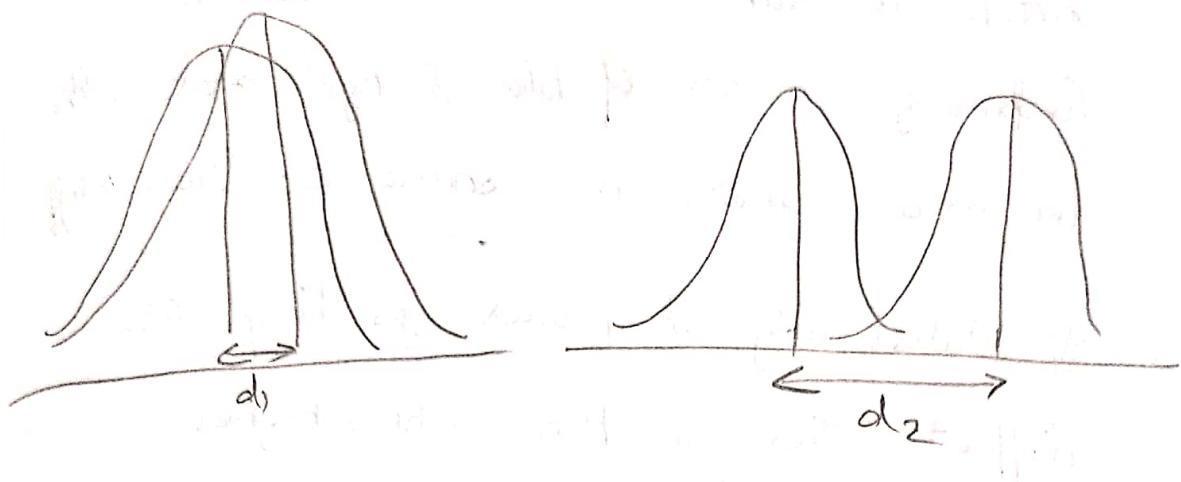
### Exercise -

<u>tstat</u>	<u>Pval</u>	<u>Null Hypothesis</u>
-1.57	0.143	Accept $H_0$ ( $P < 0.05$ )
-1.99	0.01	Reject $H_0$ ( $P < 0.05$ )
-2.1	0.0026	Reject $H_0$ ( $P < 0.05$ )
1.23	0.438	<del>Reject <math>H_0</math></del> ( $P > 0.05$ )

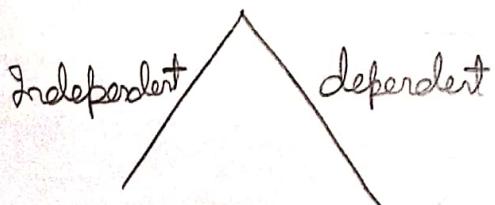
Accept  $H_0$ .



## Two Sample t - test -



- we are gonna check whether the mean of two groups are same or significantly differed.
- If the gap between the 2 means are for 95% C.I i.e if the gap is more than 1.96 then the means are almost significantly different.
- for 2 Sample t - test -



Q If the test of Blore and Hyderabad batch is taken and the mean of comparing mean of blore & hyd and with the mean score is same or signi diff.

→ If statistically it proves yes they are different, then we know who's higher.

$$H_0: \mu_{\text{Blor}} = \mu_{\text{hyd}}$$

$$H_1: \mu_{\text{Blor}} \neq \mu_{\text{hyd}}$$

If  $H_0$  accepted, the mean is same, if rejected, then we will know which batch is performing better.

This is an example of independent two sample t test.

t test - ind ( $G_1, G_2$ )

↓  
Groups.

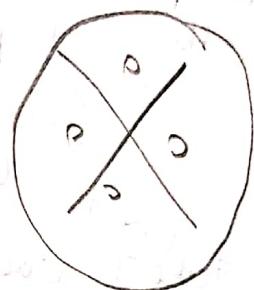
\* A test end of 4 day and found average score is 68 and next week after review, test taker and mean found out 72.

This is an example of dependent t-test

Sample t test

\* Before and after a process \*

test set  $(G_1, G_2)$   
↳ Group array



\* As we are inferring from the samples or the population, to play it safe, we never say we never say accept the  $H_0$  we always use we failed to reject null hypothesis.

IMP → we either reject or fail to reject the null hypothesis.

\* WE REJECT THE NULL

OR

\* WE FAIL TO REJECT THE NULL

HYPOTHESIS.

In Reality.

		<u><math>H_0</math>: True</u>	<u><math>H_0</math>: False</u>
Reject $H_0$ .	Type I - error $\alpha$ - error (wrong decision)	$1 - \alpha$	(Correct decision)
Accept $H_0$ .	(Correct decision) $1 - \beta$ (power of test)	Type II error $\beta$ error (wrong decision due to many things) ↓ Costly error.	

Type I error - missing out a best

Candidate due to the interview process not being good. ~~too~~

Type II error - A worst / bad candidate gets selected due to bad interview process. (Costly) (risky).

Type I error - when we reject the null hypothesis when it is true

do not reject

Type II error - when we ~~accept~~ the null hypothesis when it is false.

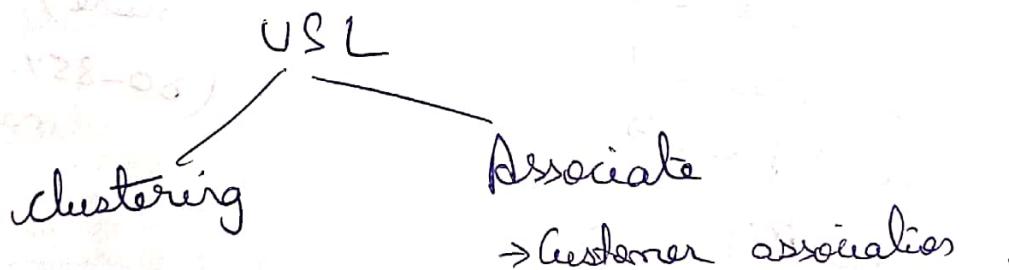
e.g'

- 1) Really missing out on a genuine good person and breaking up with her (Type I error)
- 2) Accepting a really bad person and continuing the relationship (Type II error)  
(very bad)

---

\* Supervised learning model - in this, the outcome variable is already present in the data.

\* If there is no outcome data, all are independent variables, in such data we cannot perform supervised modelling that's why we do unsupervised learning.



- 
- \* If they reject the null hypothesis they can be correct or make Type I error.
- \* If they do not reject the null hypothesis they can be correct or make a Type II error.

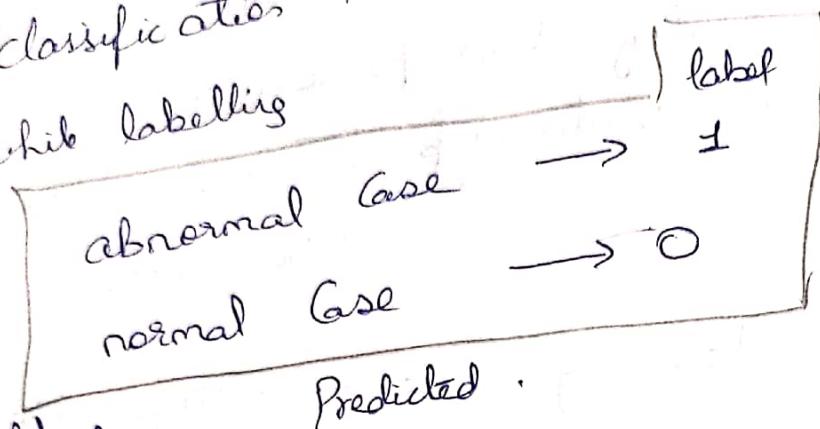
Healthy  $\rightarrow 400$   
disease  $\rightarrow 100$

Type I  
True but rejected

Type II  
False but accepted

Classification model as 2 outcomes.

while labelling



model  
specificity

		Predicted	
		0	1
Actual	Healthy	$0$	$20$
	Disease	$400$	$380$

This is people  
having disease  
is actual but  
model predicts

healthy  
(very dangerous)

$(1 - \beta)$

$60$

$\rightarrow$  model

$\frac{60}{100} = 60\%$  (very imprecise score)

not happy

(80-85% max)

Specificity - at what rate the model is predicting healthy as healthy

Sensitivity - at what rate the model predicts disease as disease (1st priority)

## ROC - Receiver Operating Characteristics

*	0	1	
0	TN	FP $\rightarrow \alpha$	True Negative False Positive
1	FN	TP	False negative True Positive
B			

$TP \rightarrow$  Sensitivity - Power of test.

$TN \rightarrow$  Specificity.

$$\text{Accuracy} = \frac{(1-\alpha) + (1-\beta)}{\text{Total records.}}$$

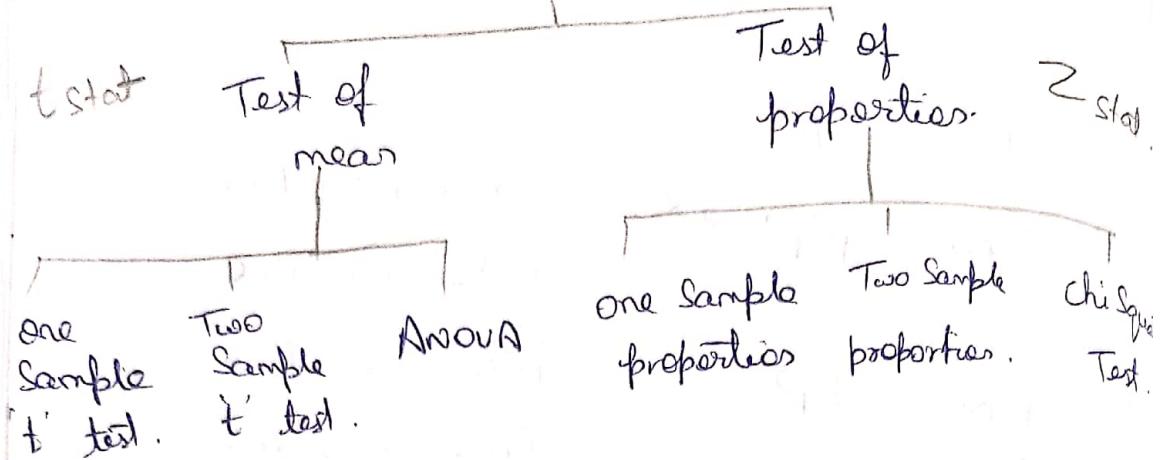
\* Even though the overall accuracy of the model can be good, that doesn't mean that the model is good, the sensitivity is the of the model is very geo. important.

Here  $\rightarrow$  only  $\frac{60}{100} = 60\%$  Sensitivity is less.

$$\text{but overall accuracy} = \frac{380 + 60}{500} = \frac{440}{500} = 88\%.$$

Still this model is not performing well for this data.

# Statistical test



$$t_{\text{stat}} = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \cdot \frac{s}{\sqrt{n}}$$

$$H_0: \mu > 81$$

$$H_1: \mu \leq 81$$

$t_{\text{stat}}$  is referred to as  $Z$  test if the population S.D is given to us.

$$t_{\text{stat}} = \frac{\bar{y} - \mu}{\frac{s}{\sqrt{n}}}$$

population S.D.

One Sample  $\rightarrow H_0: \mu > 81$

$t$  test  $\quad H_1: \mu \leq 81$

Two Sample  $\rightarrow H_0: \mu_{\text{ind}} = \mu_{\text{sig}}$   
 $t$  test  $\quad H_1: \mu_{\text{ind}} \neq \mu_{\text{sig}}$

$$P < 0.05$$

reject  $H_0$ . There is significant difference in salary.

One Sample  $\rightarrow$   $H_a: P_{\text{diabetic}} \geq 25\%$ ,  
properties test  
 $H_0: P_{\text{db}} \leq 25\%$ .

(1D) proportion  
Two Sample  $\rightarrow$  proportion of girls in chennai is equal to  
properties test  
test

(2D)  $\chi^2$  square  $\rightarrow$  prop of girls in chennai is equal to  
test  
prop " " in hyderabad.

(multi dim)

Two Sample t-test  $\rightarrow$  Same Sample

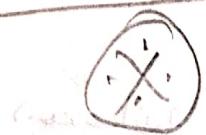
Paired (before & after Same Sample)  $n = \text{Same}$

Unpaired (Independent) (t-test - ind.)

---

\* If data strictly follows normal distribution  
go with parametric test.

Otherwise (go with non parametric test.)



test for normality

Shapiro test -

From Scipy.stats import ~~shapiro~~ shapiro.

$H_0$  : Data = normal

$H_1$  : Data  $\neq$  normal.

If e.g -  $P = 0.837$

fails to reject null hypothesis as  $P > 0.05$

$\therefore$  The data is normal.

we can continue with parametric test.

\* Even if one group is not normal, then we have to go with non parametric

non parametric Two Sample  
test  $\rightarrow$  't' test.

paired

wilcoxon

unpaired  
(Independent)

manwhitneyu

from Scipy.stats import manwhitneyu  
wilcoxon.

\* for Two Sample t test,

first with Shapiro test check for normality. If the data is normally distributed, go with parametric test depending on what kind of sample -

if data = normally dist.

paired t test  $\rightarrow$  t test - rel.

unpaired t test  $\rightarrow$  t test - ind

Even if one of the group is not following normal distribution, go with non-parametric

test.

If data not normally dist.

paired t test  $\rightarrow$  wilcoxon.

unpaired t test  $\rightarrow$  manwhitney u.

Thabu

( $\alpha, \beta$ ) loc - test

Ses - 2

\* Start with assumptions  $H_0$  is True.

for Two Sample test -

Conditions to be met -

- 1) Random variable (random Sample)
- 2) Data is normal. (Shapiro test)
- 3) Test of Variance  $\sigma_{g_1}^2 = \sigma_{g_2}^2$

Only if the above three conditions are met will go for parametric test.

parametric test

Independent

$ttest\_ind(g_1, g_2)$

Dependent / Paired

$ttest\_1samp(g_1 - g_2, 0)$   
 $ttest\_rel$

if any one condition fails, we do non-parametric

non-parametric test

Independent  
(mannwhitneyu)

Dependent (upward)  
(wilcoxon)

\* If any one of the 2 conditions fails.

→ If data is normal then to check the variance, then do 'Levene test'.

→ If data is not normal, just to gain extra insight, we do 'Bartlett test'.

$H_0$ :

Test of Variance

Data is normal.

Data is not normal.

'Levene'

Bartlett

$$H_0: \sigma_{g_1}^2 = \sigma_{g_2}^2$$

$$H_a: \sigma_{g_1}^2 \neq \sigma_{g_2}^2$$

← null hypothesis &  
alt hyp for  
variance test.

(HR.txt), names = ['HR', 'age', '']

[file.csv, index\_col = 0]

### Test of Normality -

$H_0$ : Data = Normal

$H_a$ : Data  $\neq$  Normal.

### Test of Variance - Levene / Bartlett

$H_0$ :  $\sigma_{g_1}^2 = \sigma_{g_2}^2$

$H_a$ :  $\sigma_{g_1}^2 \neq \sigma_{g_2}^2$

### Test of mean - ttest, ind / mannwhitneyu

$H_0$ :  $\mu_{g_1} = \mu_{g_2}$

$H_a$ :  $\mu_{g_1} \neq \mu_{g_2}$

Average age of male = 36.65  
" female = 37.32



acceptable range should be  
either 1.96.

$$37.32 - 36.65 = 0.67.$$

### Sessions 3 -

\* We only reject null hypothesis or  
fail to reject null hypothesis.

95%	$P < 0.05$	Reject null hypothesis
99%	$P < 0.01$	Reject $H_0$ .

## Proportion Test

$$Z_{\text{data}} = \frac{P_1 - P_2}{\sqrt{P_{\text{pooled}} * (1 - P_{\text{pooled}}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$P_{\text{pooled}} = \frac{x_1 + x_2}{n_1 + n_2}$$

$H_0$  - assess the proportion of female employees attrition is equal to male employee attrition.

$$H_0: P_{\text{FA}} = P_{\text{MA}}$$

$$H_a: P_{\text{FA}} \neq P_{\text{MA}}$$

more females are leaving or more male are leaving or irrespective of gender they are equal.

\* Attrition is not dependent on gender, its independent.

↳ Statement.

\* The moment there are no numerical columns in the test, t test is ruled out.

$P_1$  - Prob of male emp. leaving Comp.

$P_2$  - " " female " " "

$$P_1 = \frac{150}{882}$$

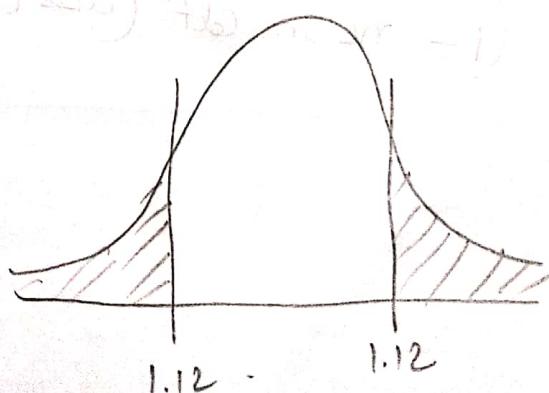
$$P_2 = \frac{87}{588}$$

$$P_{\text{pool}} = \frac{(150 + 87)}{(588 + 882)}$$

$$= 0.16 \text{ chance of leaving the job}$$

$$Z_{\text{data}} = \frac{P_1 - P_2}{\sqrt{P_{\text{pool}} * (1 - P_{\text{pool}}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$Z_{\text{data}} = 1.129$$



As this is done manually we are use the normal distributions to get an approximate p-value. later on we have built in function to calc the same

$$(1 - \text{norm.cdf}(z))^* 2 \text{ for both sides.}$$

\* So for this we got p value = 0.258.

→ so, for this we fail to reject the  $H_0$ .

Thus  $H_0$  holds good for this.  
which means, <sup>There is no significant</sup> the proportion of female att is equal to prop of male att.

even if we get negative statistic value or negative  $z_{\text{data}}$  —

$$(1 - \text{norm.cdf}(\text{abs}(z_{\text{data}})))^* 2$$

\* Verify Gender is having any influence in headache.

\* In python ↑.

Built in functions →

CHI - SQUARE TEST for Goodness

of Fit.

$$\chi^2_{\text{data}} = \sum \frac{(O-E)^2}{E}$$

$$\text{expected freq} = \frac{(\text{row total})(\text{column total})}{\text{grand total}}$$

\* Properties test for more than 2 groups is done by chi square test.

E.g (data)  $\rightarrow$  Observed (Contd.)

	Healthy	Mild	Severe	Total
Male Sample 1	410	340	250	1000
Female Sample 2	95	85	70	250
Total	505	425	320	1250

Females -

$$P_{FH} = \frac{95}{250} = 0.38 = 38\%$$

$$P_{FM} = \frac{85}{250} = 0.34 = 34\%$$

$$P_{FS} = \frac{70}{250} = 0.28 = 28\%$$

Males -

$$P_{MH} = \frac{410}{1000} = 0.41 = 41\%$$

$$P_{MM} = \frac{340}{1000} = 0.34 = 34\%$$

$$P_{MS} = \frac{250}{1000} = 0.25 = 25\%$$

Consider females -

\* In order for the null hypothesis to be true ~~what~~ is the properties required is called the expected Count.

If there is a large difference in the expected Count and original observed Count, then it indicates that the alternate hypothesis holds good.

as the difference is greater. If the original observed Count itself is needed not much correct i.e. observed Count is almost equal to expected

the  $(O-E)$  becomes 0, then  $H_0$  ~~star~~ holds good, there is no significant diff.

Expected Count

$\frac{1000 \times 505}{1250} =$ 404	$\frac{1000 \times 425}{1250} =$ 340	$\frac{1000 \times 320}{1250} =$ 256	1000
$\frac{250 \times 425}{1250} =$ 101	$\frac{250 \times 320}{1250} =$ 64	64	250
505	425	320	1250

GT

\* we are not modifying the original dof

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

$$= \frac{(404-410)^2}{404} + \frac{(340-340)^2}{340} + \frac{(256-250)^2}{256} + \\ \frac{(101-95)^2}{404} + \frac{(85-85)^2}{85} + \frac{(64-70)^2}{64}$$

$$\chi^2 = 1.1486.$$

\* The p-value for Chi-Square dist is diff from Z + dist.

Built in functions -

The data should be in cross tab or D/R, for 2-groups -

chi2 - Contingency (D/R/CT)

from Scipy.stats import chi2, chi2\_contingency

n row = degrees of freedom, dof = n of group  
= 3 - 1  
= 2.

chiValue, p-value, dof, EC =  $\chi^2$  - Contingency (Cross Tab)

$\downarrow$                      $\downarrow$   
 degrees            expected  
 of freedom            Count

Q Type of headache:

	female	male
aura	$\frac{1593}{3545} 44.9\%$	$\frac{117}{607} 19.2\%$
mixed	$\frac{291}{3545} 8.2\%$	$\frac{166}{607} 27.3\%$
No aura	$\frac{1661}{3545} 46.8\%$	$\frac{324}{607} 53.3\%$

This is the gross tab we have is further -

	female	male
aura	1593	117
mixed	291	166
No aura	<u>1661</u>	<u>324</u>
	<u>3545</u>	<u>607</u>

\* Males are highly sensitive with mixed type of migraine. as we can see females are only 8.2% whereas male 27.3%.

\* whereas females are less sensitive to aura type female - 44.9% male - 19.2%.

Questions -

- Heart data - (2)  
disease  
state ?
- 1) Gender influences heart state ? (2)  
↳ Two Sample properties test  
(2)  
with heart disease
- 2) Effect of exercise (2)  
↳ Two Sample properties test.  
(continuous)  
with cholesterol test
- 3) Any sig difference in cholesterol test (2)  
exercised people  
↳ we have to mean of exercise and  
non exercise people  
↳ Two Sample test.

HR data -

- $\chi^2$
- 1) Any particular dept having high att rate
- 2) Any particular education.  $\chi^2$
- 3) check the prob of m/f across all  
depts same?

df [ 'Education' ]. Value - Counts ()  $\rightarrow$  one dimensional.

chi square  $\rightarrow$  for one dimensional.

$\rightarrow$  It is going to tell whether the proportion of no's is same across all the group or different.

$\rightarrow$  If we put chi square in a single column suppose department, then it will tell us the properties of department is same or not.

54  
50  
56  
52  
54

It will tell the department is properties or not.

Gender	Male	Female
Age	18-25	26-35
Education	10th	12th

Gender	Male	Female
Age	18-25	26-35
Education	10th	12th

chi square.

Gender	Male	Female
Age	18-25	26-35
Education	10th	12th

One variable but multi Categ we are doing with samples.

## 1-Samp proportion test

25% emp percentage of Indians women  
are diabetic in India.

$$H_0 = 0.25$$

$$H_a \neq 0.25$$

we perform Sampling  $n = 500$

$$\text{Diab.} = 200$$

$$H = 300$$

Fun-

from stats models. stats. Properties  
import proportions - Z test

Proportions - Z test (200, 500, 0.25)

↓      ↓      ↓  
diabetic    Total    percentage  
sample      sample    expected

\* This can be used for two sample prop test also.

Contd -  $(x_1, x_2)$

Obs -  $(n_1, n_2)$

one Sample -

1) Proportion - Z test (Actual Count, tot Samples, claim on population)

Two Sample -

2) Proportion - Z test (Obs, Count, Total sample, to check claim made on population, yes, alternative, exercise)

→ Sample, proportion test.

3) Chi 2 - Contingency (CT)

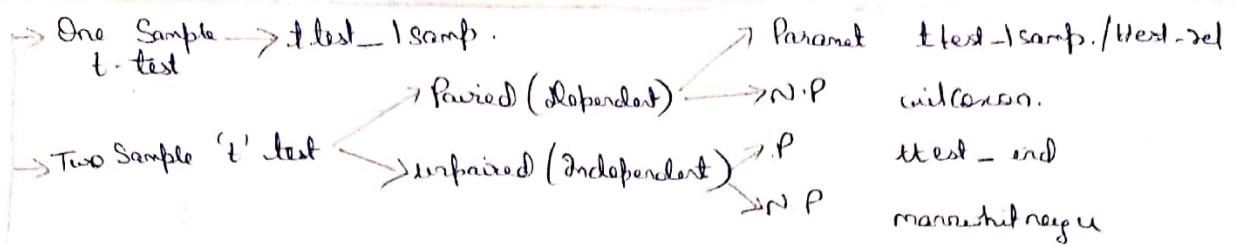
4) Chi Square (Value Counts single Category)  
(Sampling).

→ To check whether the sample is a representation of the populations.

Proportion - Z test (Affected, Total)

## STATISTICAL TEST

Test of mean.



- Two Samples (groups)  
w.r.t one categorical variable  $\rightarrow$  f - One way

Two groups w.r.t two (or) more categorical variable  $\rightarrow$  2-way ANOVA O/S library

Test of proportion.

One Sample Z test with single category (Proportions - Z test)

One Sample Z test with multiple category (chi-square)

Two Sample Z test (Proportions - Z test)

Two or more groups (chi<sup>2</sup> - Contingency)

# Day 4

## Sessions 1 -

### ANALYSIS OF VARIANCE

Sample ages for Groups A, B, C.

<u>Group A</u>	<u>Group B</u>	<u>Group C</u>	
30	25	25	
40	30	30	
50	50	40	
60	55	45	
mean 45	40	35	Global mean = 40

Sample ages for Groups D, E, F

<u>Group D</u>	<u>Group E</u>	<u>Group F</u>	
43	37	34	
45	40	35	
45	40	35	
47	43	36	
mean - 45	40	35	

$$G.M = 40$$

# Mean Square Treatment (MSTR)

↳ Between Sample Variability.

$$MSTR = \frac{\sum n_i (\bar{x}_i - \bar{\bar{x}})^2}{k-1} \text{ or } \frac{SSTR}{dof}$$

Between Sample Variability

$n_i$  = Sample no.

$k-1$  = dof

for g - ABC -

$$H_0: \text{Mage}(A) = \text{Mage}(B) = \text{Mage}(C)$$

$$H_a: \text{Mage}(A) \neq \text{Mage}(B) \neq \text{Mage}(C)$$

for group - DEF

$$H_0: \text{Mage}(D) = \text{Mage}(E) = \text{Mage}(F)$$

$$H_a: \text{Mage}(D) \neq \text{Mage}(E) \neq \text{Mage}(F)$$

$$MSTR = \frac{4(45-40)^2 + 4(40-40)^2 + 4(35-40)^2}{3-1}$$

$$= \frac{25+0+25}{2} \times \frac{4(25)+0+4(25)}{2}$$

$$MSTR = 100$$

Between Sample Variability

Within Sample Variability

Sample Variability

$$F_{\text{stat}} = \frac{MSTR}{MSE} = \frac{\frac{SSTR}{df_1}}{\frac{SSE}{df_2}}$$

Variane of  $A = 166.66 \left(\frac{1}{n}\right)$

$B = 216.66 \left(\sigma_B^2\right)$

$C = 83.33 \left(\sigma_C^2\right)$

\* while using  $\text{np.var}(A, \text{ddof}=1)$  always use ddof.

because it calculates for

$$\frac{\sum_{i=1}^n (v_i - \bar{x})^2}{n} \quad \text{whereas we need } \frac{\sum_{i=1}^n (v_i - \bar{x})^2}{n-1}.$$

$$MSE = \frac{SSE}{df_2} = \frac{\sum (n_i - 1) s_i^2}{n_t - k} \rightarrow \begin{array}{l} 3 \text{ group} \times 4 \text{ sample} \\ \text{per group} \\ = 12 \text{ samples.} \end{array}$$

$(n_i - 1) = 3$

$s_i = \sigma_A^2 + \sigma_B^2 + \sigma_C^2$

$$= \frac{3 \times \sigma_A^2 + 3 \sigma_B^2 + 3 \sigma_C^2}{12 - 3} \rightarrow \text{groups!}$$

total  
Samples

$$\begin{aligned}
 &= \frac{(3 \times 166.66) + (3 \times 216.66) + (3 \times 83.33)}{9} \\
 &= \frac{499.98 + 649.98 + 249.99}{9}
 \end{aligned}$$

$$MSE = 155.55$$

$$\begin{aligned}
 f_{\text{stat}} &= \frac{MSTR}{MSB} \\
 &= \frac{100}{155.55}
 \end{aligned}$$

$$f_{\text{stat}} = 0.6428$$

$$P_{\text{value}} = 0.548$$

$$P > 0.05$$

$\therefore$  We fail to reject the null hypothesis.

$$\therefore \text{Mag}(A) = \text{Mag}(B) = \text{Mag}(C)$$

Python functions -

\* The data should be in form of dataframe

model -

mod = ols('age ~ group', data = df).fit()

↑

always left side  
should be continuous

A -  
A -  
A -  
B -  
B -  
B -  
C -  
C -  
C -

for Anova table -

which numerical var a/c to  
which cat variable  
 $\text{aov\_table} = \text{Sm. Stats. aov} - \text{lm}(\text{mod, typ} = 1)$

print(aov\_table)

mod = ols('Continuous ~ Categorical', df).fit() display  
for anova  
mod

for group - D, E, F

MSTR is same

$$\text{MSTR} = 100.$$

Variance of  $D = 2.666 (\sigma_D^2)$   
 $E = 6.00 (\sigma_E^2)$   
 $F = 0.666 (\sigma_F^2)$

$$\text{MSE} = \frac{\text{SSE}}{\text{df}_2} = \frac{\sum (n_i - 1) s_i^2}{n_t - K}$$

$$= \frac{3 \times \sigma_D^2 + 3 \times \sigma_E^2 + 3 \times \sigma_F^2}{12 - 3}$$

$$= \frac{(3 \times 2.66) + (3 \times 6.00) + (3 \times 0.66)}{9}$$

$$= \frac{7.98 + 18 + 1.98}{9}$$

$$\text{MSE} = 3.10$$

$$F_{\text{STAT}} = \frac{MSTR}{MSE}$$

$$= \frac{100}{3.10}$$

$$F_{\text{STAT}} = \underline{\underline{32.25}}$$

Calculating on python -

$$\text{Pvalue} = 0.00008$$

$$\therefore \text{Pvalue} < 0.05$$

∴ we ~~do~~ Reject the null hypothesis

Here -

$$\text{Mage}(A) \neq \text{Mage}(B) \neq \text{Mage}(C)$$

\* First we should check whether the anova test is passing or not, if the mean values for all the groups look same, then we cannot predict the group based on age. So, in order to predict group based on age, then there should be difference in mean values of the groups.

\* from Scipy.stats import f\_oway  
The columns must be grouped by first  
To get complete detail, use anova\_lm.

$$F_{\text{stat}} = \frac{\text{Between Sample Variation}}{\text{within Sample Variation}} = \frac{MS_{\text{B}}}{MS_{\text{W}}}$$

- \* The first group (A, B, C) mean is same whereas, the second group (D, E, F) mean is different
- \* The  $H_0$  for group 1 got accepted whereas the  $H_0$  for group 2 got rejected because of the more spread of data in group 1, the variance for that is more, whereas the spread of data in group 2 is less as compared to 1.

Inference -

B

- group (ABC)
- group (DEF)

A

Here we can see that the spread of data of group 1 is higher than the spread of data of group 2

\* Since, the datapoints of group 1 have too much of spread is there, there is a high risk of these two groups overlapping each other, they are not significantly different.

but as the spread of data for group 2 is very less, then the chance of them overlapping is very rare, hence

they are ~~there is no~~ significant diff.

\* If the between variability increases for group 1, then, it can ~~become~~ be concluded that  $M_D + M_B \neq M_E$ .

Inference -

The between Sample Variability for both (ABC) & (DEF) share the same value, but within Sample Variability for (DEF) was very compact, because of that their means are far apart.

$$\left( \frac{MSR}{MSE} \right)$$

$$F_{\text{stat}} = \frac{100}{155} = 0.6428$$

$$F_{\text{stat}} = \frac{100}{3.11} = 32.14$$

→ Both  $g_1$  &  $g_2$  are sharing same numerator which means the dist b/w groups is same (between Sample Variability is same)

→ But within Sample Variability for  $g_1$  is more for  $g_2$  than  $g_2$  (compact)

→ If the ABC mean should be sign. diff then either the denominator should shrink or the numerator should increase more.

Check avg. Sal vs Dept.  
 check average salaries in all 3 departments  
 for this we use anova test

	df	df	Sum Sq	mean Sq	F	PR(>F)	Pr value
group.	2.0	SSTR	200	MSTR 100.00	Fstat 0.64285		0.5483
Residual	9.0	SSE	1400.0	MSE 155.556		NAN	NAN

## Two Way Anova

\* Simultaneously  
 categories with  
 columns, then it is  
 called two way  
 Anova.

when we check two  
 mean of a continuous

called two way

alone cannot predict salary

Salary ~ Dept.

Interactions

Salary ~ Joblevel

if put together is  
 more any change

$H_0$  : There is no interactions

$H_a$  : There is an interactions

$$Sal = \beta_0 + \beta_1 Joblevel$$

$$Sal = \beta_0 + \beta_1 Department \rightarrow \begin{array}{l} \text{This is ruled out} \\ \text{as we cannot pred Sales,} \\ \text{by dept} \end{array}$$

$$Sal = \beta_0 + \beta_1 Joblevel + \beta_2 (Department \times Joblevel)$$

$$mod = ols \left( \text{Salary} \sim Job\_level : \text{Dept} \right), \text{fit}$$

$$Anova\_tab = \text{Anova\_lm} (mod, \text{typ}=1)$$

In two way Anova, here, the  $H_0$ : There is no interaction b/w the two columns is getting rejected, here we deduce that, yes, both put together influence the predictions of the Salary.

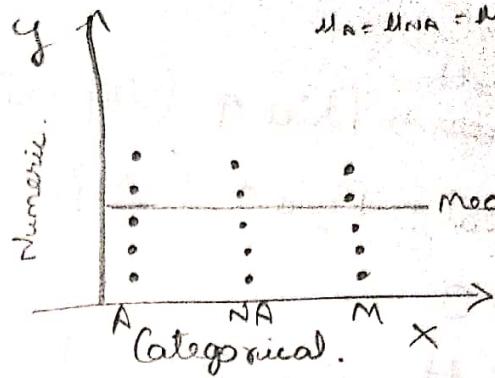
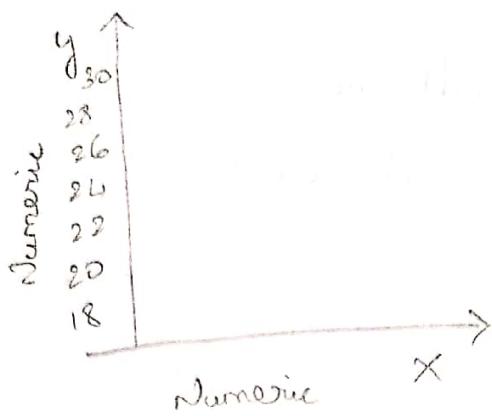
So, If we use departments and job-level as an interaction, it has a better chance at predicting it (Salary)

Correctly.

Only the product of department and Sales has an influence on predicting Salary.

$$Sal = \beta_0 + \beta_1 (Joblevel \times Dept)$$

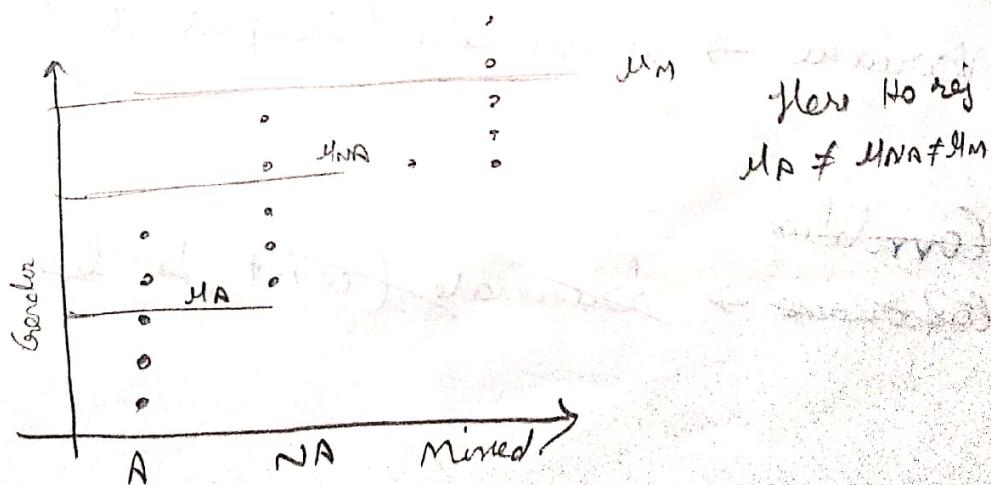
# Regressives -



Aura  
NAura  
Mixed

# 9 Here the mean of all the three samples of all the three samples are same, so,  $H_0$  holds good  $\mu_P = \mu_{NA} = \mu_M$ , so this is a useless variable as predicting.

In order to reject the null hypothesis to, the mean of the 3 categories used should be different, only then will be able to use it predictors of the dependent var. (age)



Here  $H_0$  rej  
 $\mu_P \neq \mu_{NA} \neq \mu_M$

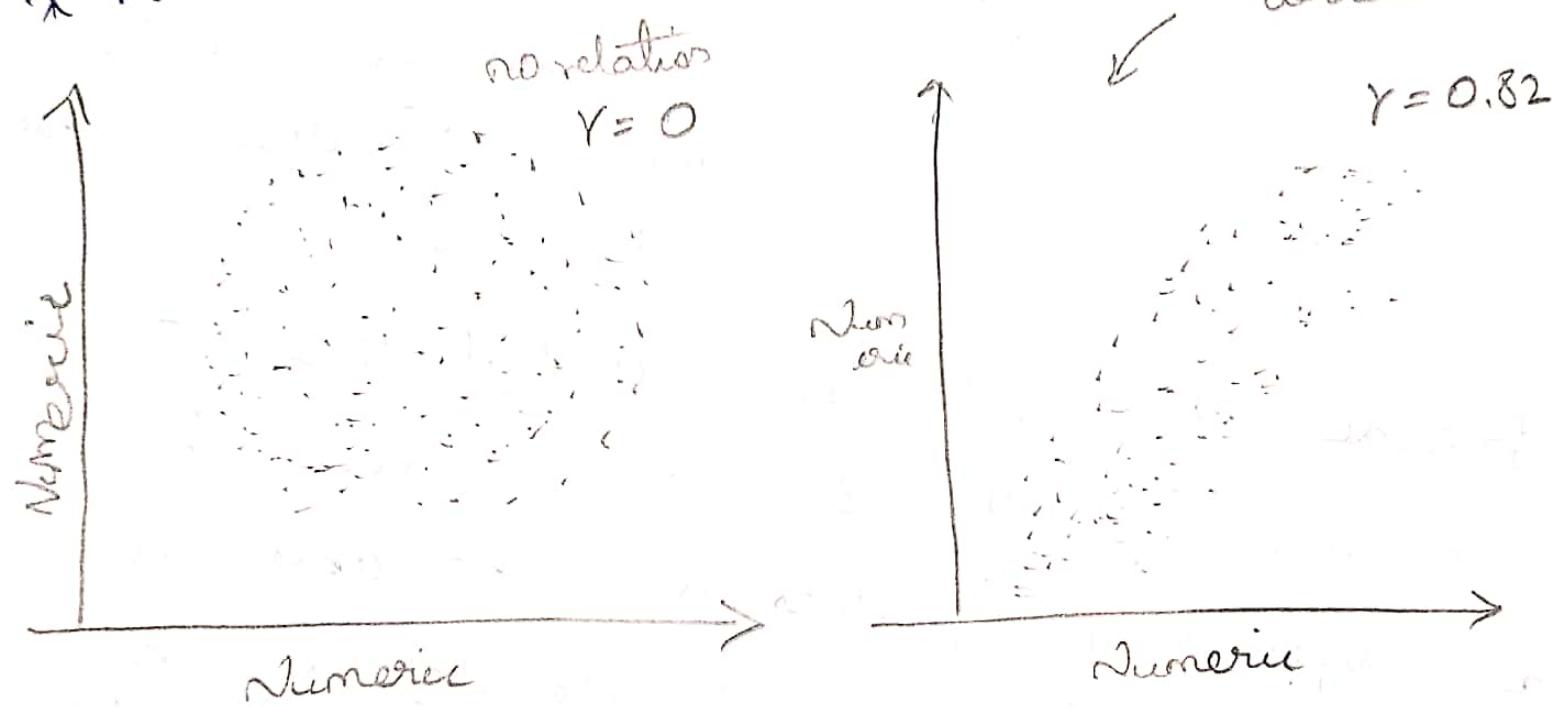
$$\text{Price} = f(L, B, L \times B)$$

~~~~~ \curvearrowleft Interactions (Area)

- \* ~~length~~  $\rightarrow$  Price  $\uparrow$  (Separately influence)
- \* ~~Breadth~~  $\rightarrow$  Price  $\uparrow$  (Separately influence)
- \*  $L \times B$  means area, but include, either length or breadth one, or use the  $L \times B$  for predictions, otherwise the data will be redundant.

\* Numeric VS Numeric -

Positively  
Correlated

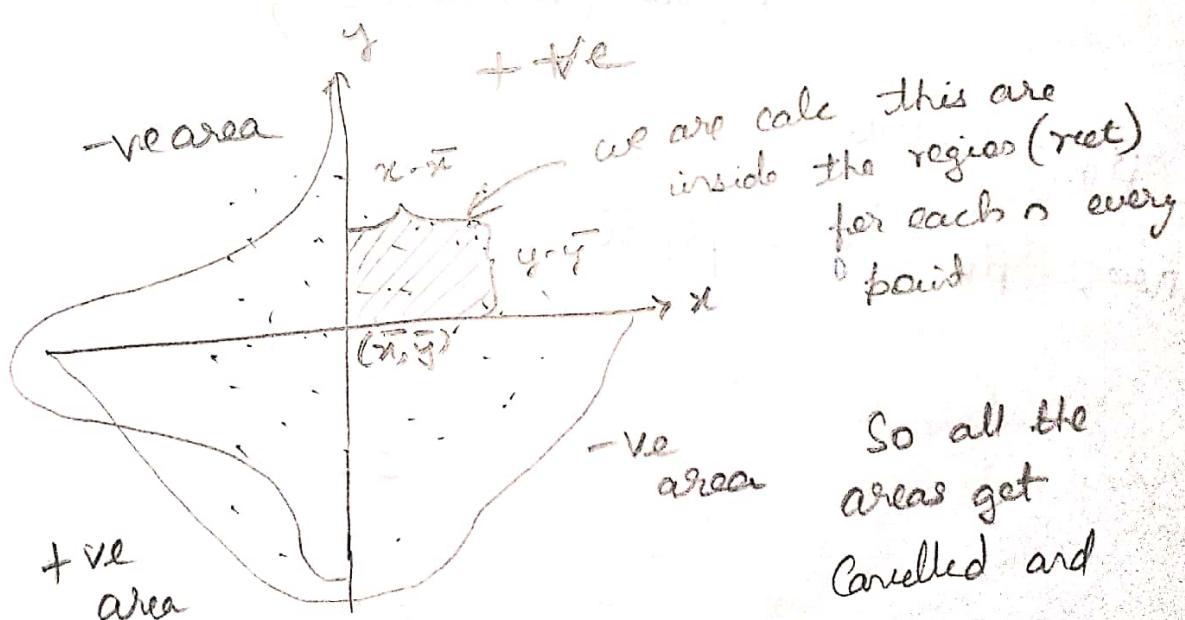
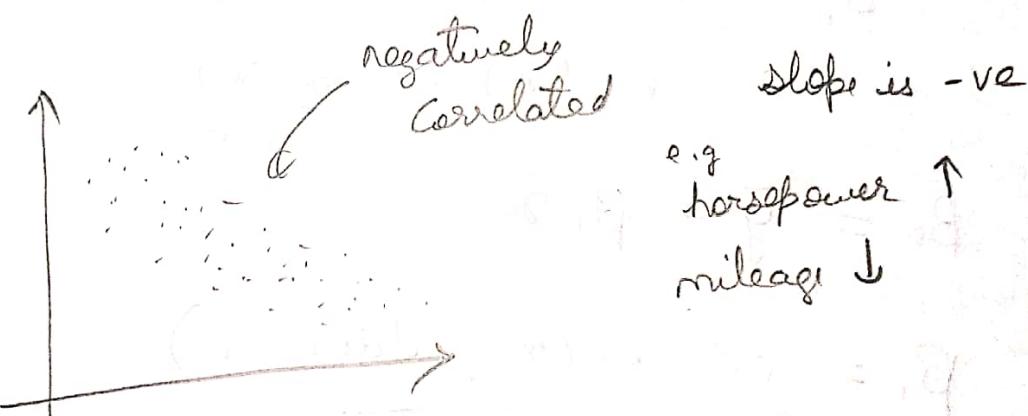
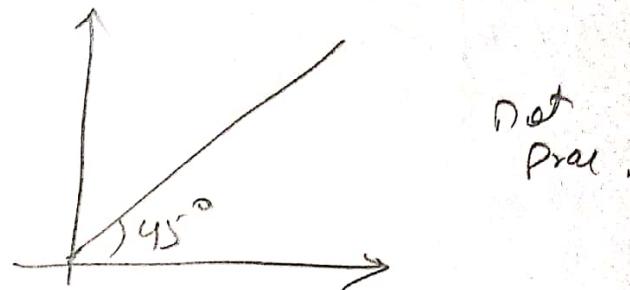


Variance  $\rightarrow$  univariate (respect to one variable)

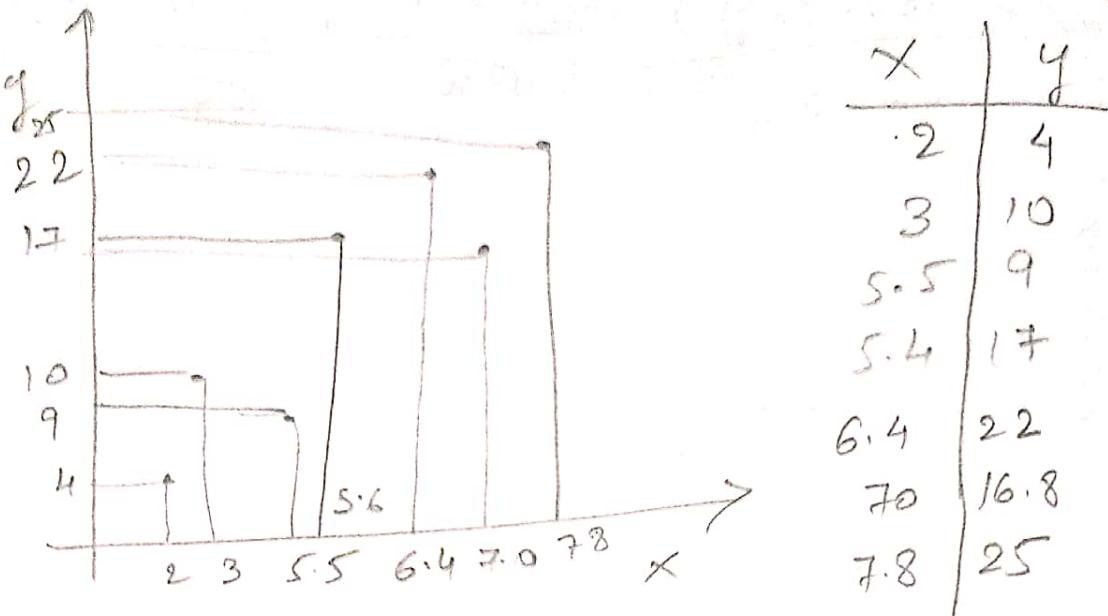
$$\gamma_{xy} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{x - \bar{x}}{\sigma_x} * \frac{y - \bar{y}}{\sigma_y} \quad (z_{\text{score}} \text{ form})$$

| $\frac{x}{3}$ | $\frac{y}{3}$ |
|---------------|---------------|
| 3             | 3             |
| 5             | 5             |
| 8             | 8             |

$\gamma_{xy} = 1$  Self Correlated.



So all the areas get Cancelled and  $\gamma = 0$ .



$$\beta_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{\frac{(x - \bar{x})(y - \bar{y})}{n}}{\frac{(x - \bar{x})^2}{n}}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_1 = \frac{n \cdot \text{Cov}(x, y, \text{ddof}=1)}{n \cdot \text{Var}(x, \text{ddof}=1)}$$

$$\beta_1 = 3 \cdot 10^6$$

now,  $\beta_0 = \bar{y} - \beta_1 \bar{x}$

| x   | y    | $\bar{y}$ |
|-----|------|-----------|
| 2   | 4    | 4.57      |
| 3   | 10   | 7.60      |
| 5.5 | 9    | 15.44     |
| 5.4 | 17   |           |
| 6.4 | 22   |           |
| 7.0 | 16.8 |           |
| 7.8 | 25   |           |

for modelling -

$x = 2$  dimension

$y = 2D / 1$  Dim.

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
|       |       |     |

$$x = DF [x_1' + x_2']$$

$$y = DF [y']$$

model = ols('y ~ x1 + x2', DF). fit()

y\_Pred = model.predict(x)

model.params

$\beta_0$

$\beta_1$

$\beta_2$

\* Statistical  
analysis  
of  
variables

$x = DF [x_1, x_2]$  or  $DF[DF[$   
 $x_1, x_2], y]$

$y = DF [y']$

model = Linear Reg()

model.fit(x, y)

y\_Pred = model.predict(x)

model.Coeff -

\* Not  
possible to  
do stat analyses

# Hypothesis for model building

$$H_0 : \beta_1 = 0$$

$$H_a : \beta_1 \neq 0$$

$$R^2 = \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

denom: large to get  
good  $R^2$

e.g -  $Age = \beta_0 + \beta_1 (\text{exp})$

$$\beta_0 = \bar{y} - \beta_1 \bar{x} \quad \beta_1 = \frac{\text{cor}(x, y)}{\text{var}(x)}$$

if the spread of the variable to be predicted is less, then it will lead to the poor  $R^2$  value. So, the denominator should be large.

$$R^2 = \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

- \* The average difference is the error.
- \* The average error made by the model is RMSE
- \* The smaller the RMSE, better the model.
- \* Larger the  $R^2$ , better the model.

|                                           | $R^2$ | RMSE |
|-------------------------------------------|-------|------|
| $y = \beta_0 + \beta_1 Rad.$              | 51%   | 4.25 |
| $y = \beta_0 + \beta_1 TV$                | 61%   | 3.24 |
| $y = \beta_0 + \beta_1 TV + \beta_2 Rad.$ | 89%   | 1.66 |

This is the best model.

OLS → to get summary

practical

deployment is done using Sckit learn.

RMSE built in func -

from sklearn import metrics

metrics. mean\_squared\_error