

07
 WEEK
 DAY-1:-

STATISTICS:-

25th Week Day 173-192

Saturday

ML

$y \rightarrow$ dependent variable
 $x \rightarrow$ independent var.

Prediction Model.

Regression

$$y = f(x)$$

Nature of y here is
 Continuous values

Classification

$$y = f(x)$$

Here the nature of y is 0 or 1
 (categorical).
 $y = \beta_0 + \beta_1 x$

$$n \times p$$

no my cols

newly fields

entire features

instances attributes

samples headers.

variables.

$$y = \beta_0 + \beta_1 x$$

$\beta_0 \rightarrow$ Slope.

Intercept

$\beta_1, \beta_0 \rightarrow$ Model Parameters.

$x \rightarrow$ Predictor variable.

Entropy \rightarrow Measure of uncertainty

\rightarrow Historical data also known as Empirical data.

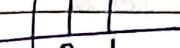
$y \rightarrow$ Dependent var | Target col | Outcome var

NOTE! There can be more than 1 outcome var.

To predict y_1 and y_2 then you need to build one model to predict y_1 and 2nd model to predict y_2 .

\rightarrow For a CATEGORICAL VAR you can make use BAR CHART or PIE CHART only. And you can only count the categorical variables.

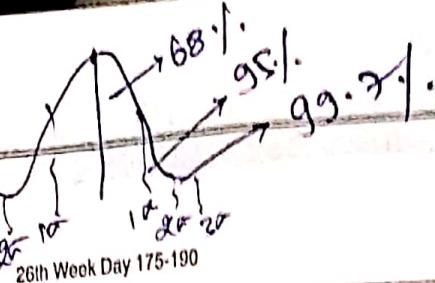
The diff b/w BAR and HISTOGRAM is BAR takes discrete values on X-axis and Histogram has contd values on X-axis.

BAR PLOT \rightarrow  Histogram \rightarrow 

JUNE

24

Monday



26th Week Day 175-190

M	T	W	T	F	S	S
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

06

OUTLIER: Any value greater than 3σ ($3 \times \text{std. dev.}$)

→ You need to know whether an outlier is an outlier or an extreme value.

OUTLIER

- ① Remove
- ② Replace with median
- ③ Resampling

EXTREME

① Log TRANSFORM.

→ Extreme values are important
→ OUTLIERS are NOISE and unwanted

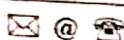
→ Exponential transformation is applicable when all of your data is in the range of 10^0 . This is also not a good representation of data.

Inferential Stats

→ When your sample data infers about the entire population.

Descriptive Stats

→ Talks about the columns in the data only.



TRUE MEAN: Your mean and my mean remains the same.

Mean of the population

importance in our life is greater (u).
and now mean \rightarrow God knows.

2019 WK	M	T	W	T	F	S	S
07	1	2	3	4	5	6	7
27	8	9	10	11	12	13	14
28	15	16	17	18	19	20	21
29	22	23	24	25	26	27	28
30	29	30	31				

JUNE

25

26th Week Day 176-183

Tuesday

$N \rightarrow$ population size $n \rightarrow$ Sample size

$\mu \rightarrow$ population mean (true mean) $\bar{x} \text{ or } \bar{y} \rightarrow$ Sample mean. \bar{y}

$\sigma \rightarrow$ Standard deviation $s \rightarrow$ Standard devi. of sample

Note To perform inferential stats the minimum sample must be $n > 30$ and has to be highly random.

Hypothesis Testing

① Test of Mean ② Test of Proportion

done for numerical data done for Categorical data.
(T-test) (Z-test)

Continuous

Categorical

Univariate { Histogram
Box plot

Bar plot

Pie chart.

Bivariate
Analysis

Scatter
plot

df. conc.).

Quantiles for

$$Q_1 = \frac{1}{2} (\text{Sample}) = \frac{93}{2} = 46.5 \approx 47^{\text{th}} \text{ elem}$$

$$Q_3 = \frac{3}{4} (\text{sample}) = \frac{93}{4} = 23.25 \approx 24^{\text{th}} \text{ elem.}$$

$$Q_5 = \frac{3}{4} (\text{sample}) = \frac{3}{4} \times 93 = 69.75 \approx 70^{\text{th}} \text{ elem.}$$

JUNE

26

Wednesday

26th Week Day 177-188

Shallow COPY :-

$$A = [1, 2, 3]$$

$$B = A$$

$$B[1] = 20$$

$$\text{print}(B) = [1, 20, 3]$$

$$\text{print}(A) = [1, 20, 3]$$

changes made in B are reflected in A.

Probability :-

Used Majorly \leftarrow (1) Classical Probability :-

(2) Empirical Probability (Historical Data).

(3) Subjective Probability.

Probability Distribution :-

→ when you have only 2 outcomes then it is a binomial distribution.

→ when you have many numbers of outcomes then its a Poisson distribution.

Calculating Probability in Binomial Distribution

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

$$\rightarrow p(x=1) = {}^5 C_1 (0.05)^1 (0.95)^4$$

Binomial
Poisson
normal

M	T	W	T	F	S	S	Wk Jun 2019
1	2						22
3	4	5	6	7	8	9	23
10	11	12	13	14	15	16	24
17	18	19	20	21	22	23	25
24	25	26	27	28	29	30	26

06

27	1	2	3	4	5	6	7
28	8	9	10	11	12	13	14
29	15	16	17	18	19	20	21
30	22	23	24	25	26	27	28
31	29	30	31				

07

26th Week Day 178-187

Thursday

Formula for Poisson Dist.

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

$\lambda = \text{No. of events in a given time}$

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

where, $f(x)$ is used to represent a probability density fun.

μ is any value $-\infty < \mu < \infty$.

$$\sigma = 2.8, \quad \mu = 3.14$$

Central Limit Theorem
 → Majority of the data lies closer to the mean
 b/w - 1σ and + 1σ zone.

Standard Deviation for Normal Dist -

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

Binomial Dist -

$$SD = \sqrt{np(1-p)}$$

Poisson Dist - $SD = \sqrt{\lambda}$

Random Variable:

Categorical

Discrete

Continuous

Yes/No

No. of rooms in a house

Price of House

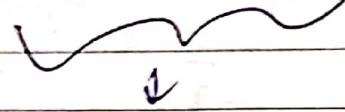
True/False

Class in Titanic

Rate of Speed

0/1

Def/ND.



You make me
of Poisons DST

You get multiple
OUTCOMES.

You make me

of NORMAL DST

You make use of

BINOMIAL

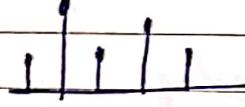
DST



You get only
2 OUTCOMES

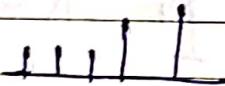
The values are cont.

Binomial



Pmf (prob. mass fnⁿ)

Poisson



Pmf (prob. mass fnⁿ)

Normal



pdf (prob. density fnⁿ)

NOTE: In a normal dist. if I wish to calculate $P(X = 45\text{cm})$ is almost ≈ 0 . Because in a Normal Dist you find the Density. (Area under the curve).

~~for one single point AREA under the curve is zero~~

~~In normal dist you need two values to get the probability.~~

Jul 2019	W	M	T	W	T	F	S
		1	2	3	4	5	6
		7	8	9	10	11	12
		13	14	15	16	17	18
		19	20	21	22	23	24
		25	26	27	28	29	30
07		31					

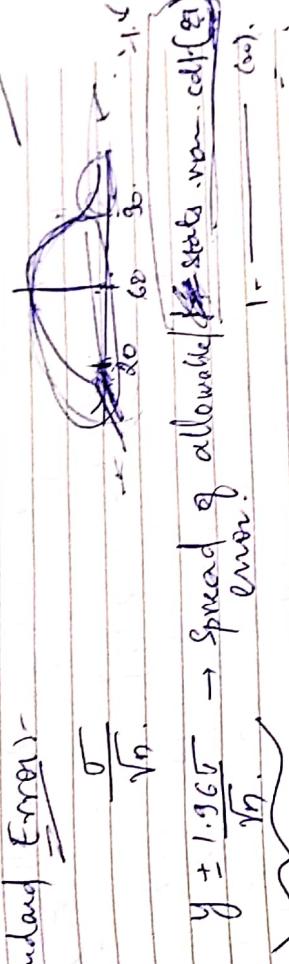
29

Saturday

26th Week Day 180-185

Standard Normal Dist:-

Scaling of data and making it ~~scale~~ unit less will make two changes mean = 0 and $SD = 1$.

Standard Error:-

$$Y \pm 1.96\sigma \rightarrow \text{Spread of allowed state norm. conf. (95%).}$$

Convers 95% of data.

Given Company claims delivery time < 36m.

Sample avg say emy:

$$n = 50$$

$$\bar{x} = 2.8 \text{ hrs}$$

$$S = 0.6 \text{ hrs} \quad \text{Now we need to infer the true mean } (\mu)$$

which will be in a range which confidence interval required is 95%.

Confidence interval required for 95% confidence the range will be

$$\boxed{\bar{x} - \frac{1.96\sigma}{\sqrt{n}}, \bar{x} + \frac{1.96\sigma}{\sqrt{n}}} \rightarrow \text{Confidence interval Range}$$

$$\text{for 99% } \left[\bar{x} - \frac{2.58\sigma}{\sqrt{n}}, \bar{x} + \frac{2.58\sigma}{\sqrt{n}} \right] \\ \text{for 98% } \rightarrow 2.33 \quad \text{for 90% } \rightarrow 1.65.$$

01

Monday 27th March Day 132 183

1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				07



- 1.96 28.96 1.96

$$\therefore \text{Confidence Interval} = \left[2.8 - \frac{1.9(0.6)}{\sqrt{50}}, 2.8 + \frac{1.9(0.6)}{\sqrt{50}} \right]$$

$$= [2.63, 2.96]$$

We observed from 2.63 to 2.96 forms 95% confidence interval for H. In other words we are 95% confident the avg delivery time lies b/w 2.63 and 2.96. Since all the values are less than 28. Thus we have a strong evidence that the claims are correct.

→ There is 5% error this is called α -error.

→ To get 1% error we compute confidence for 99.)

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

$$= [2.58, 2.01]$$

To conclude which interval is true for our population then we need to perform HYPOTHESIS TESTING.

TEST OF MEAN (or) TEST OF PROPORTION.

Aug 2019 Wk	M	T	W	T	F	S
	31	1	2	3	4	t-test
32	5	6	7	8	9	10 11
33	12	13	14	15	16	17 18
34	19	20	21	22	23	24 25
35	26	27	28	29	30	31

JUL

02

27th Week Day 183-182

Tuesday

$$n = 26, \bar{y} = 573, s = 124$$

$$\text{Confidence with 95%: } \left[\frac{573 - 1.9(124)}{\sqrt{26}}, \frac{573 + 1.9(124)}{\sqrt{26}} \right] \\ = [532.49, 613.5].$$

When $n \geq 30$ according to Statistical theory

$s \approx \sigma$
where $s \rightarrow \text{Standard deviation of sample}$
 $\sigma \rightarrow \text{Standard deviation of population}$.

→ more or -
↳ t-distribution is same as normal distribution
↳ It is used during hypothesis testing.

$$\text{For normal dist } Z = \frac{x - \mu}{\sigma}$$

$$\text{For t-dist } Z = \frac{x - \bar{x}}{\frac{s}{\sqrt{n}}}.$$

HYPOTHESIS TESTING

$H_0 \rightarrow$ Alternate Hypothesis (or) Research Hypothesis (or) H_{A}

In H_0 the claim is made on the population.

@ $H_0 \rightarrow$ Null Hypothesis. (a) No Difference by pop.

$H_0: \mu_0 > 520$ (Right Tail)
for $H_0: \mu_0 < 320$ $H_0: \mu_0 \leq 520$ (Left Tail).

$H_0: \mu_0 \geq 320$.

JULY THIS is One Sample t-test.

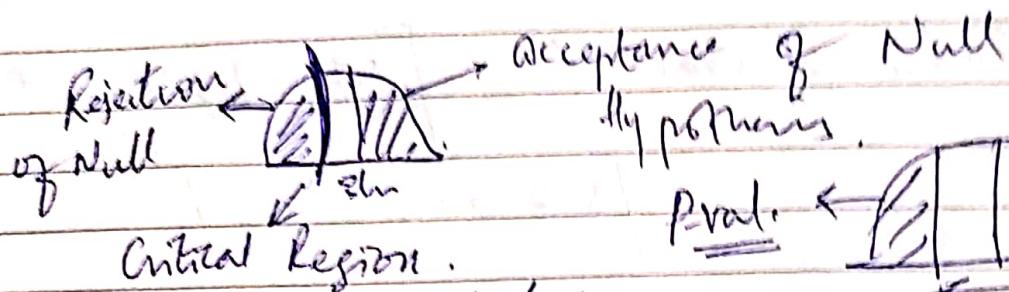
03

Wednesday

27th Week Day 184-181

M	T	W	F	S	S	Wk Jul 2016
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				
						07

$$H_a : \mu_a < 3 \text{ hr} \quad H_0 : \mu_a \geq 3 \text{ hr}$$

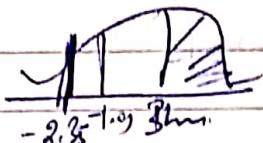


Critical Region is calculated as:

$$t_{\text{stat}} = \frac{\bar{y} - \mu}{s/\sqrt{n}} = \frac{2.8 - 3}{0.6/\sqrt{50}} = -2.35$$

RULE for REJECTING NULL HYPOTHESIS

	Should be 95% ggg.	Should be 99%
t_{stat}	> 1.96	> 2.58



P-value	< 0.05	< 0.01

P-value \rightarrow Probability of null hypothesis being TRUE.

One Sample t-test :-

One Sample \rightarrow Population sample.

t-test & Test of Mean using t-distribution.

JULY

REGISTRATION NO.	NAME	ROLL NO.	SECTION
08	S. S.	1 2 3 4	
	S. S.	5 6 7 8 9 10 11	
	S. S.	12 13 14 15 16 17 18	
	S. S.	19 20 21 22 23 24 25	
	S. S.	26 27 28 29 30 31	

04

27th Week Day 165-180

Thursday

Two Sample t-test :-

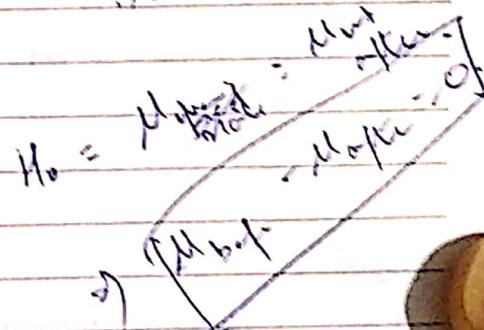
$$t_{\text{stat}} = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

For 2 Sample t-test we have = or \neq criteria.
Only for 1 Sample t-test we have $>$, $<$, \geq , \leq .

$$\Rightarrow n=20, \bar{y}=2.28, s=0.65 \quad H_0: \mu < 3.84, \\ H_a: \mu > 3.84$$

$$t_{\text{stat}} = \frac{2.28 - 3}{0.65 / \sqrt{20}} = \frac{-0.22}{0.65 / \sqrt{20}} = -0.22$$

$$\Rightarrow t_{\text{stat}} = -2.85$$



Type 1 and Type 2 Errors:-

Type 1 error is α error

Type 2 error is β error.

DECISION	NULL		HYPOTHESIS
	$H_0: \text{TRUE}$	$H_0: \text{FALSE}$	
Reject H_0	(Wrong decision) Type 1 error α -error	Correct Decision	$1 - \alpha$.
Accept H_0	Correct Decision $1 - \beta$	Type 2 error β -error	

Type 2 error is a costly error.

JULY

05

Friday

27th Week Day 100-170

P value is the Probability of Null Hypothesis being TRUE.

M	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

$$H_0 \rightarrow \text{Good} \quad D \rightarrow 100 \quad H_0: \text{TRUE} \quad H_0: \text{FALSE}$$

		0	1
Reject H_0	Neckly	$(1-\alpha) TN$	αFP
	$\frac{1}{100}$	380	20
Accept H_0	Disease	βFN	TP
	$\frac{1}{100}$	50	$(1-\beta)$

$\rightarrow (1-\beta) \rightarrow$ sensitivity of the Model.

$$\therefore \text{Sensitivity of the model} = \frac{60}{100} = 60\%$$

$$1-\alpha \rightarrow \text{Specificity of the model} = \frac{380}{400} \approx 95\%$$

$$1-\beta \rightarrow \text{Sensitivity of the model.}$$

Notes: First importance is given to sensitivity.

Sensitivity is also referred to as TRUE POSITIVE RATE
Specificity is also referred to as TRUE NEGATIVE RATE

Stat Test

t test

Test of Mean

z test

Test of Proportion

One Sample

Two Sample

ANOVA

One Sample

Two Sample

Chi-Square Test

't' test

't' test

'prop'

'prop'

t test, samp (S, μ)

t test

'prop'

Chi-Square Test

$t_{stat} = \bar{y} - \mu$

$\bar{y}_1 - \bar{y}_2$

Unpaired

$\frac{s}{\sqrt{n}}$

$\sqrt{\frac{s_1^2 + s_2^2}{n_1 + n_2}}$

(Independent)

Paired

(Independent)

$t_{stat}, \text{samp } (g_1 - g_2, 0)$

ANOVA \rightarrow More than 2 groups.

If Normal \rightarrow Parametric

Vice versa.

Aug 2019 Wk	M	T	W	T	F	S
31	1	2	3	4		
32	5	6	7	8	9	10 11
33	12	13	14	15	16	17 18
34	19	20	21	22	23	24 25
35	26	27	28	29	30	31

08

JULY

06

27th Week Day 187-178

Saturday

~~(1)~~ Sometimes t-test referred as z-test when population Standard Deviation is given to us.

ttest_1samp ($g_1, g_2, 0$), ttest_ind () \rightarrow parametric.

These are followed only if your data is normal or not.

SHAPIRO TEST :-

Test for normality. (Check if your data is normal or not).
from scipy.stats import ttest_1samp, ttest_ind, Shapiro.

\rightarrow Shapiro will return the p value.

Shapiro(age).

if $p = 0.837$.

H_0 : data - Normal

H_a : data \neq Normal.

Null hyp is accept the H_0 .

Now if the data is normal then you can proceed with parametric test.

\rightarrow If it is not normal, then we will go with non-parametric test.

Sunday 07

Two Sample t-test

Paired

Unpaired

✉ @ ✉

Wilcoxon

ttest_1samp ($g_1, g_2, 0$)

ttest_ind

MannWhitneyU

Non-parametric

Parametric

Parametric

Non-Parametric

JULY 08 If all 3 cond's are satisfied, then you go for parametric test.

Monday

28th Week Day 189-176

Two Sample t-test:

- ① Random Variable (Random Samples)
- ② Data is Normal (Normality Test)
- ③ Test of Variance ($\sigma_1^2 = \sigma_2^2$)

Go for parametric test.

Independent

t-test-ind (g_1, g_2)

Dependent

t-test-remp ($g_1, g_2, 0$)

→ If any one cond' fails, you do non-parametric

Non-Parametric

Independent
Mann-Whitney

Dependent

Wilcoxon

Test of Variance.

Data is Normal

(Levene)

Data is Not Normal

(Bartlett)

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

$$H_a: \neq$$

AUG 2019 Wk	M	T	W	T	F	S	S
	31	1	2	3	4		
	32	5	6	7	8	9	10
	33	12	13	14	15	16	17
	34	19	20	21	22	23	24
	35	26	27	28	29	30	31

Wilcoxon & Mann - are used very often.
→ used 90% of the times.

JULY

ANOVA → Analysis of Variance

09

28th Week Day 190-175

Tuesday

Test of Normality
(Shapiro)

$$H_0: \text{Data} \sim \text{Normal}$$

$$H_a: \text{Data} \neq \text{Normal}.$$

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 / Mannwhitney)

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

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

Test of PROPORTION :-

$$Z_{\text{stat}} = \frac{P_1 - P_2}{\sqrt{\text{Pooled} \cdot (1-\text{Pooled}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

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

CHI SQUARE TEST :-

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

O → observed count E → expected count.

$$\text{Expected frequency} = \frac{(\text{Row total})(\text{Column total})}{\text{Grand total.}}$$

from statsmodels.stats.proportion import proportion_2test.

counts = np.array([10, 87])

obs = np.array([882, 588])

proportions.2test(counts, obs).

10

Wednesday

28th Week Day 191-174

1	2	3	4	5	6	7	27
8	9	10	11	12	13	14	28
15	16	17	18	19	20	21	29
22	23	24	25	26	27	28	30
29	30	31					31

07

Example Data:-

	H	M	S
Male	410	360	250
Female	95	85	70

	H	M	S	} Observed Count
Male	410	360	250	
Female	95	85	70	

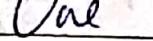
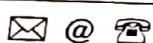
Expected

Count =

404	340	256	} Obtained using the formula for expected count.
101	85	64	

Chi².contingencyChi².contingency

Chi-square.



One Sample Prop. Test :-

$$H_0: 0.25$$

(25% of middle age women being diabetic)

$$H_a \neq 0.25$$

Aug 2019 Wk	M	T	W	T	F	S	S
31		1	2	3	4		
32	5	6	7	8	9	10	11
33	12	13	14	15	16	17	18
34	19	20	21	22	23	24	25
35	26	27	28	29	30	31	

Test of Proportion at a Glance

JULY
11

20th Week Day 192, #73
actual

Thursday

- (*) For one sample proportion test proportions - χ^2 test (~~total count, total samples~~, claim proportion on population).
- (*) Two sample proportion test. proportions - χ^2 test (counts, obs) $\xrightarrow{\text{actual count}} [n, n_f]$ $\xrightarrow{\text{[n, n_f]}}$
- (*) \rightarrow 2 sample proportion test. chis-contingency (crosstab)

- (*) chisquare (value - counts of single categorical variable) used mostly in sampling process.

H HD Ex vs HD.

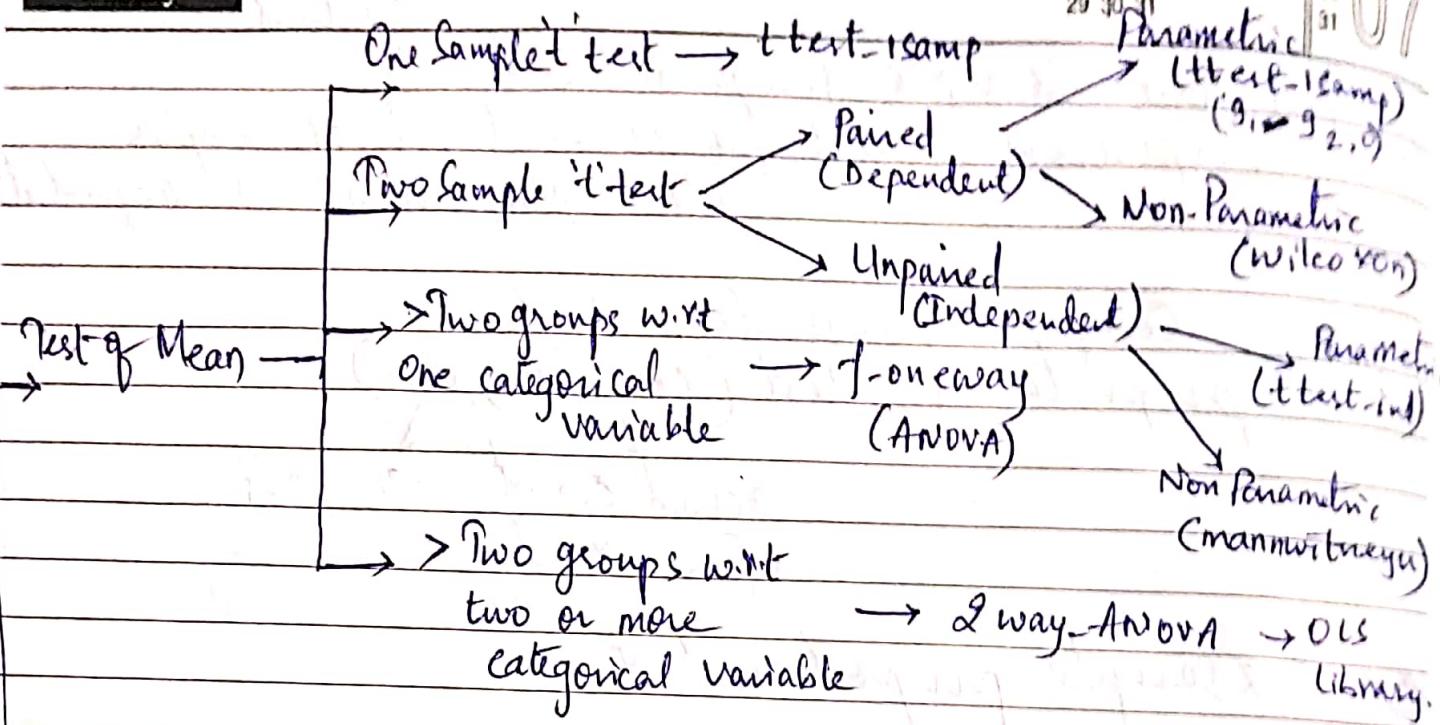
E	127	54	$n_1 \rightarrow$ total exercise
	23	66	$n_2 \rightarrow$ total not exercise
NE	n_1	n_2	$\xrightarrow{\text{total HD}}$
	n_1	n_2	\rightarrow HD

proportions - χ^2 test $\left([54, 66], [127, 89] \right)$.

✉ @ ☎ Here $H_0: P_{\text{HD}}(\text{E}) = P_{\text{HD}}(\text{NE})$

1	2	3	4	5	6	7	27
8	9	10	11	12	13	14	28
15	16	17	18	19	20	21	29
22	23	24	25	28	27	26	30
29	30	31					31

07



STATISTICAL TEST

→ One Samp Z test with (Proportions - z-test)

→ One Samp - Z test with (chisquare)

→ Two Sample z test (Proportions - ztest) M F

Test of
Proportion

✉ @ ☎

→ Two groups (chi²-contingency) A B C

→ Test of
Proportion

31	1	2	3	4
32	5	6	7	8
33	9	10	11	
34	12	13	14	15
35	16	17	18	
	19	20	21	22
	23	24	25	
	26	27	28	29
	30	31		

08

~~a Glance~~~~Everything at~~

13

28th Week Day 194-171

Saturday

(Q1) Whether heart disease is having any impact on the cholesterol.

here $H_0: \mu_{CH} = \mu_{CNH}$ (Test of Mean)

$H_a: \mu_{CH} \neq \mu_{CNH}$

① Two sample t-test.

② It is a independent hence unpaired

③ If not parametric go

(Q2) Effect of Vessel on cholesterol level.

Here $H_0: \mu_{C(0)} = \mu_{C(1)} = \mu_{C(2)} = \mu_{C(3)}$.

$H_a: \mu_{C(0)} \neq \mu_{C(1)} \neq \mu_{C(2)} \neq \mu_{C(3)}$.

Clearly we have cholesterol as numerical
and vessels has 4 Categorical

Sunday 14

Therefore you have 4 groups.

Hence we use f-one way ANOVA.

JULY ONE NUMERIC with 2 CATEGORICAL

2 way ANOVA.

~~Mugwau Rotka~~

29th Week Day 196-169

	MON	TUE	WED	THU	FRI	SAT	SUN	WK	JUL	2019
more than 2 categories	1	2	3	4	5	6	7	27		
headache	8	9	10	11	12	13	14	28		
& categorical	15	16	17	18	19	20	21	29		
	22	23	24	25	26	27	28	30		
	31								07	

- * Whether age is a factor influencing type of headache.

age → numerical
headache → types (Categorical)

clear cut f- One way ANOVA

- * age vs headache vs gender

~~ANOVAS~~
is 2 way ANOVA.

- * dependency b/w drivetrain and no. air bags in car.

Drivetrain → Categorical

Air bags → Categorical.

Test of Association.

$$H_0: P_{FI} = P_{FO} = \dots = P_{FS}$$

This is Chi-square Contingency. Find the residuals and sum it to Chi-square Contingency.

- * Note: If null hyp does not hold then you conclude that there is a dependency.

V.T.P @ 20

Aug 2019 Wk | M T W T F S S

08 31 1 2 3 4
32 5 6 7 8 9 10 11
33 12 13 14 15 16 17 18
34 19 20 21 22 23 24 25
35 26 27 28 29 30 31

Maths behind Anova

JULY

16

29th Week Day 197-168

Tuesday

MSTR - Between Sample Variability :-

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

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

$$MSE = \frac{\sum (n_i - 1) s_i^2}{n_t - k} = \frac{SSE}{df_2}$$

$s_i \rightarrow$ Standard deviation. $df_2 \rightarrow$ no. of couples, no. of group

$$\text{For Variance} : \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

In python its 'n' in the denominator to correct it you need to pass np.var(A, ddof=1).

By default it is 'n-1' in R.

Covariance : (Bivariate).

$$\text{Cov}(x, y) = E(x - \bar{x})(y - \bar{y}).$$

$$\Rightarrow \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}.$$

✉ @ ☎ Cov (wt, ht) \Rightarrow 250 kg cm

Cov (wt, age) \Rightarrow 125 kg yrs

To remove the units you divide by SD.

$$r_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} \rightarrow \text{Varies from -1 to 1.}$$

JULY

17

Wednesday

29th Week Day 198-167

M	T	W	T	F	S	S	Wk Jul 2019
1	2	3	4	5	6	7	27
8	9	10	11	12	13	14	28
15	16	17	18	19	20	21	29
22	23	24	25	26	27	28	30
29	30	31					31

07

Linear Regression:

$$y = \beta_0 + \beta_1 x$$

$$\beta_0 \rightarrow \bar{y} - \beta_1 \bar{x}, \quad \beta_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

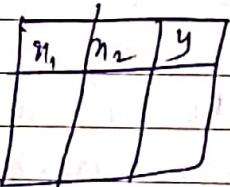
For Scikit Learn Approach:-

```
from sklearn.linear_model import LinearRegression
model = LinearRegression()
model.fit(x, y)
y = model.predict(x).
```

Steps at a Glance:-

$$(1) X = DF[['x_1', 'x_2']]$$

$$y = DF['y']$$



OLS Approach:-

$$\text{model} = \text{ols}('y ~ x_1 + x_2', DF).fit()$$

$$y_{\text{pred}} = \text{model}.predict(x)$$

model.params → will give you $\beta_0, \beta_1, \beta_2$.

Scikit approach:-

X has to be in 2 Dimensions.

$$X = DF.drop('y')$$

$$y = DF['y']$$

Aug 2019 Wk M T W T F S S

JULY

08

31	1	2	3	4			
32	5	6	7	8	9	10	11
33	12	13	14	15	16	17	18
34	19	20	21	22	23	24	25
35	26	27	28	29	30	31	

18

29th Week Day 199-180

Thursday

Framing the Hypothesis for linear Regression

$$H_0: \beta_1 = 0 \quad (\beta_1 = \text{slope})$$

$$H_a: \beta_1 \neq 0$$

$$\beta_1 = 0$$

$$\beta_1 \neq 0$$

Performance of Regression Model

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

Higher the R^2 higher the prediction rate.

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

Best Model Selection $\leftarrow R^2$ RMSE

$y = \beta_0 + \beta_1 R$	35%	4.25
$y = \beta_0 + \beta_1 TV$	61%	3.24
$\rightarrow y = \beta_0 + \beta_1 TV + \beta_2 R$	89%	1.66

Best Model Here the multivariate clearly outperforms the univariate.

