

CONTENTS

- Probability
- Permutation & Combination
- Covariance
- Pearson Correlation
- Spearman Correlation
- Hypothesis Testing
- Point Estimate
- Chi Square Test

SCHEDULE



Probability

" Probability is a measure of the likelihood of an event.

eg: 1 Tossing a fair coin

$$P(H) = 0.5, \quad P(T) = 0.5$$

eg: 2 Sholay Coin : unfair coin

$$P(H) = 1$$

eg: 3 Rolling a dice

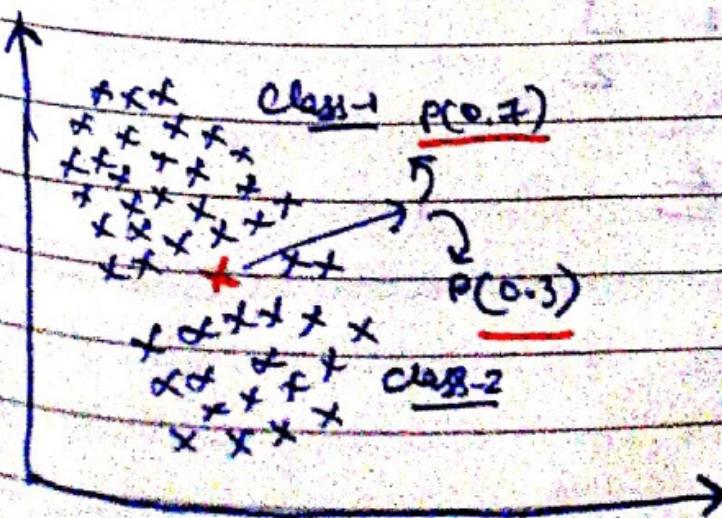
$$P(1) = \frac{1}{6}, \quad P(2) = \frac{1}{6}, \quad P(3) = \frac{1}{6}$$

* Use of Probability :-

In Machine Learning.

For eg:- like In a classification problem
 It say there are some data points belongs
 to 2 different classes. Now if we
 add a new data point then we judge
 through some algorithm that what
 will be the probability to add that
 data point in which class.

Sunday 06



Monday

09.00 ① Mutual Exclusive Events :-

10.00

11.00

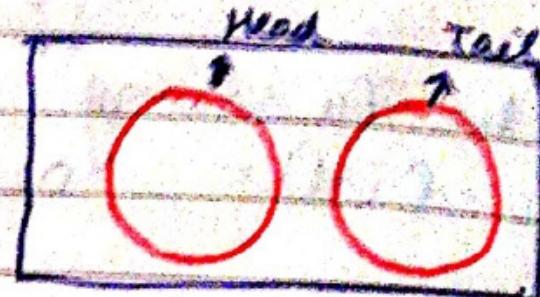
12.00

13.00

Two events are mutually exclusive if they cannot occur at same time.

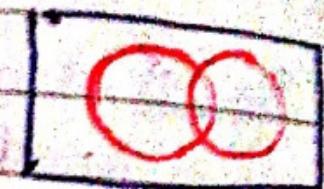
eg:-

- ① Tossing a coin
2. Rolling a dice



20.00 ②

Non-Mutual Exclusive Events



Two events can occur at the same time.

eg:- 1) Picking a card from deck of cards, two events "Heart" and "King" can be selected

Q1 what is the probability of coin landing on Heads or tail.

Hint: Use Addition Rule for mutual exclusive events.

Sol

$$P(A \text{ or } B) = P(A) + P(B)$$

$$\Rightarrow \frac{1}{2} + \frac{1}{2}$$

$$\Rightarrow 1 \approx$$

Q2

A Bag contains : 10 Red marbles, 6 Green and 3 red & green marbles.

(I) what is the prob. of choosing a marble that is red or green?

Hint: use addition rule for non-mutually exclusive event.

$$P(A \cup B)$$

Sol:

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) - P(A \text{ and } B) \\ &= \frac{10}{19} + \frac{6}{19} - P(\text{Red and Green}) \end{aligned}$$

$$R \rightarrow 10+3=13$$

$$G \rightarrow 6+3=9$$

$$\Rightarrow \frac{19}{19} = 1 \text{ E}$$

Q3 What is the prob. of choosing a 'heart' or 'Queen' from a deck of cards. [non-mutual]

$$\underline{\text{Sol}} \quad P(H) = \frac{13}{52}, \quad P(Q) = \frac{4}{52}$$

$$P(H \text{ or } Q) = P(H) + P(Q) - P(H \text{ and } Q)$$

$$\Rightarrow \frac{13}{52} + \frac{4}{52} - \frac{1}{52}$$

$$= \frac{16}{52}$$

09:00 → Multiplication Rule

10:00

11:00 ① Dependent Events :-

12:00

13:00

14:00

15:00

16:00

17:00

18:00

19:00

20:00

Two events are dependent if they affect one another.

For e.g :-

In a bag of marble there are 6 Red Balls & 4 Green Balls, $P(R) = \frac{6}{10}$
 $P(G) = \frac{4}{10}$, after removing 1 Red ball
 Prob. of green balls affect by $P(G) = \frac{4}{9}$

②

With Independent Events :-

Each and every two events is ~~one~~ independent.

for e.g.: Rolling a dice.

Q1

What is the probability of rolling "5" & "3" with a normal 6 side dice?

Hint: Use Multiplication Rule for Independent Events
 & mutually exclusive

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

$$\Rightarrow \frac{1}{6} * \frac{1}{6}$$

$$\Rightarrow \frac{1}{36}$$

Friday

09:00 → Permutation & Combination

10:00

11:00 • Permutation $\Rightarrow {}^n P_r = \frac{n!}{(n-r)!}$

12:00

13:00

14:00

15:00

• Combination $\Rightarrow {}^n C_r = \frac{n!}{r!(n-r)!}$

16:00

(unique combination)
(not repeated)

17:00

18:00

19:00

20:00

for e.g.:-

Q1 In a bag 5 chocolates are available & you have to select only 3

Sol5 4 3

$$\Rightarrow 5 \times 4 \times 3 \Rightarrow 60.$$

No. of options of selection = 60

ORQ2

$${}^n P_r \Rightarrow$$

$$[{}^5 P_3] =$$

$$\frac{15}{15-3}$$

$$\Rightarrow [60]$$

Q2 • In case of Combination: (repetition not allowed)

$${}^n C_r = {}^5 C_3 \Rightarrow$$

$$\frac{5 \times 4 \times 3}{3 \times 2 \times 1} \Rightarrow [10]$$

07 08 09
M T W T F S
SCHEDULE Imp

■ MAY ■ 2018

132-233 / Wk-19

12

Saturday

09:00

→ Covariance :

{ used in feature Selection }

10:00

11:00

12:00

13:00

14:00

15:00

16:00

17:00

18:00

19:00

20:00

$$\text{Cov}(X, Y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\text{var}_x \left[\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} \right] \Rightarrow \text{var}(x) = \frac{\sum (x_i - \bar{x})(x_i - \bar{x})}{n-1}$$

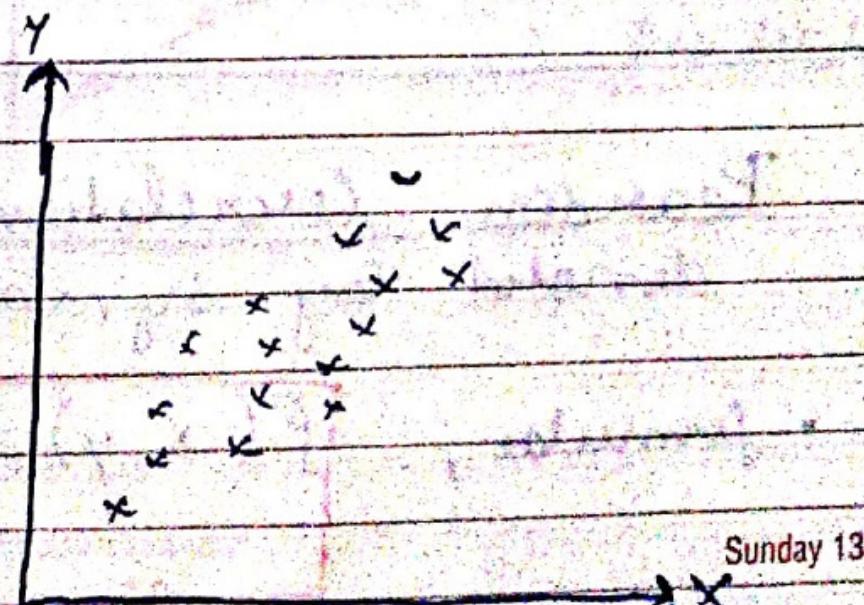
$$\Rightarrow \text{var}(x) = \text{cov}(x, x)$$

* Interview Ques : $\Rightarrow \text{cov}(x, x) = \text{var}(x)$

↳ Relationship b/w cov & var.

Eg:-

X	Y
Age	Weight
12	41
13	44
14	47
15	50
16	55
17	59
18	61



Sunday 13

\Rightarrow Age ↑ \leftrightarrow weight ↑ } observation
Age ↓ \leftrightarrow weight ↓ }

Set

$$\bar{x} = 15$$

$$\bar{y} = 51$$

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\Rightarrow 24 \rightarrow (+\text{Ve})$$

* According to final ~~Ques Ans~~ $x \propto y$

* if Cov is (-ve) then $x \propto -y$

* if, cov is '0', then $\boxed{\text{NO relation b/w } x \& y}$

→ Pearson Correlation Coefficient (P)
denoted by symbol Rho (P) [for linear data]

• formulae, $P(x, y) = \frac{\text{Cov}(x, y)}{[\sigma_x \cdot \sigma_y]}$ → we are restricting
-1 to +1 by using $[\sigma_x \cdot \sigma_y]$

* W.r.t. covariance there is no such restriction that how much +ve value it can have or -ve value it can have.
But with the help of Pearson Correlation we are trying to restrict the values between [-1 to +1]

09.00

10.00

11.00

12.00

13.00

14.00

15.00

16.00

17.00

18.00

19.00

20.00

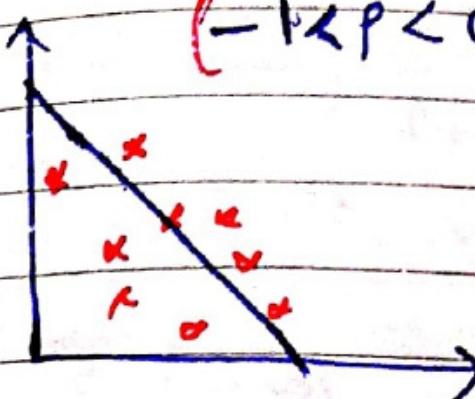
* "More the value towards +1,
more positive correlated it is."

* "more the value towards -1,
more negative correlated it is."

Eg (at wikipedia)

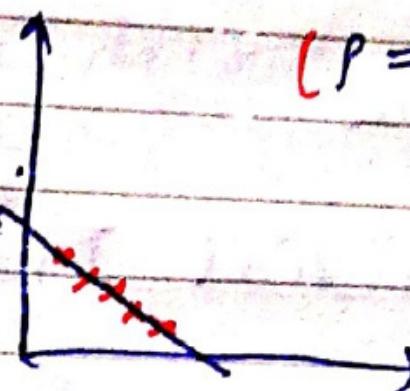
$$(-1 < \rho < 0)$$

a)



$$(\rho = -1)$$

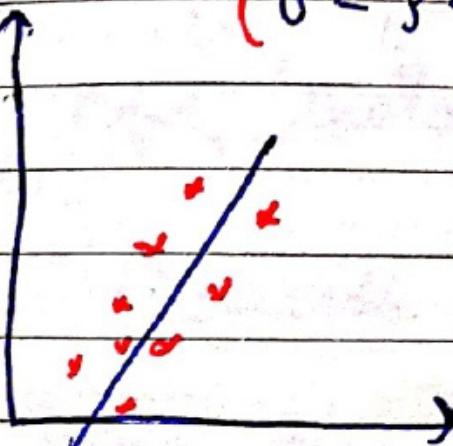
b)



$$\begin{cases} x \uparrow & y \downarrow \\ x \downarrow & y \uparrow \end{cases}$$

$$(0 < \rho < +1)$$

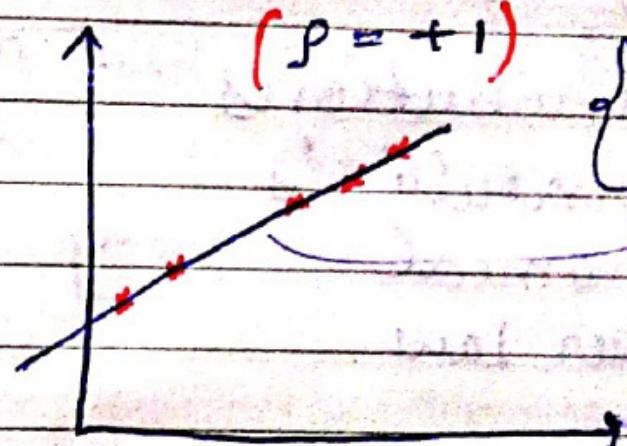
c)



$$(\rho = +1)$$

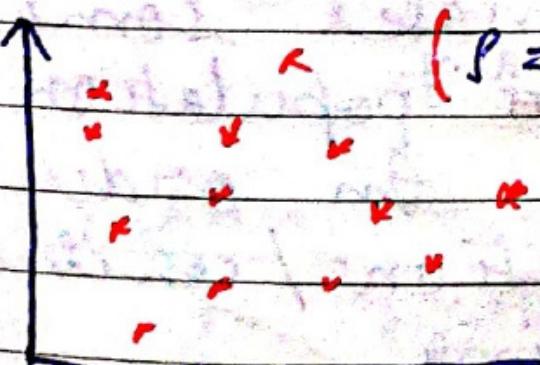
$$\begin{cases} x \uparrow & y \uparrow \\ x \downarrow & y \downarrow \end{cases}$$

d)



Best fit line

e)



$$(\rho = 0)$$

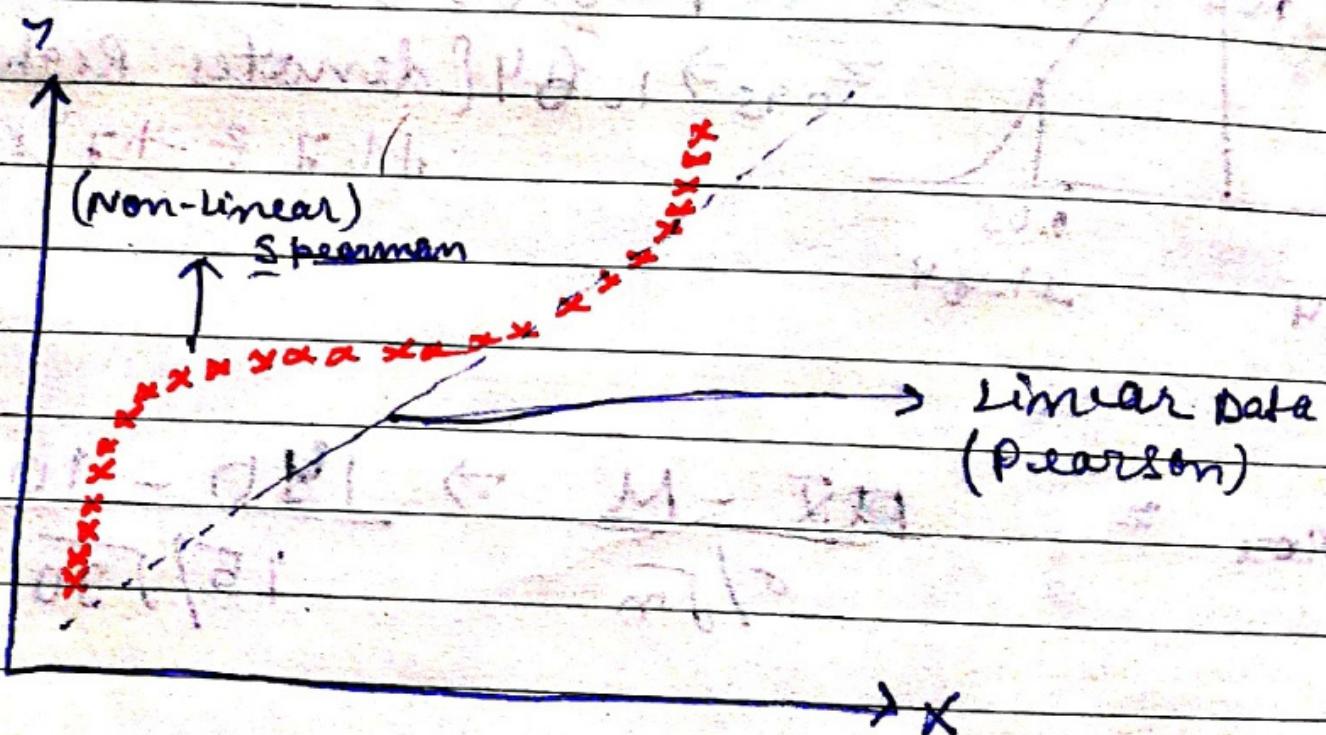
[NO Relation]

→ Spearman Rank Correlation :- {Non-Linear}

Disadvantage with Pearson Correlation is that it only holds Linear Data for entire correlation.

So for Non Linear Data points we have to use Spearman Rank Correlation. Spearman assigns Rank, formulae :- (in ascending Order) ↪

$$R_s = \frac{\text{cov}(R(x), R(y))}{\sigma(R(x)) \sigma(R(y))}$$



- Pearson Corr : 0.88 ($0 < \rho < 1$)

- Spearman Corr : 1 (Rank)

Hypothesis Testing

In order to validate assumption regarding population data through sample data we use Hypothesis Testing.

→ 2 types of Hypothesis Testing :-

• Steps of Hypothesis Testing

① Null Hypotheses :- The default one (right one)
• Eg:- Person is not a criminal until he is proven by court.

Eg 2 ② [Coin is fair or Not]

To prove it should be $P(H) = 0.5$ & $P(T) = 0.5$

• In case, of fair coin is Null Hypothesis.
(except shotay coin)

② Alternate hypothesis

- from eg-1 Person has committed a crime.
- from eg-2 coin is not fair

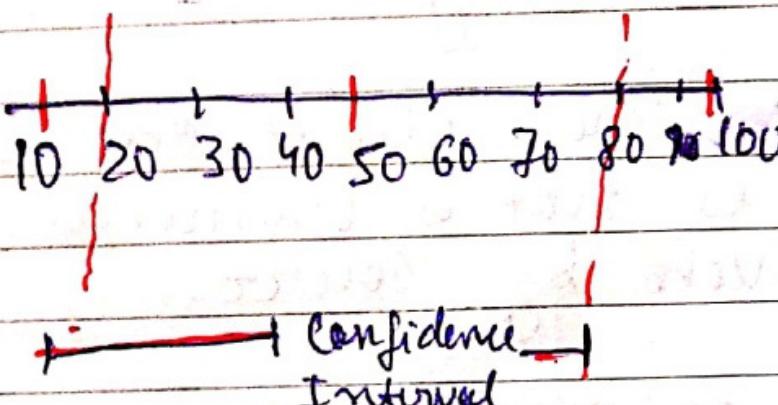
③ Perform Experiments :- Perform experiments to prove some conclusion,
tossing a coin.

- for eg. performing experiments 100 times
 - a) First time got 50 times Head : coin is fair in this
 - b) Second time got 60 times Head : Can say coin is fair
 - c) Third time got 70 times Head : " " " "

Friday

SCHEDULE
 2018-19
 5.5.18 - 18.7.19

So, conclusion is based on the range between 10 to 100. So if the range (C-I) is fixed that b/w 20 - 80 if condition lies in this interval then we can say coin is fair, outer this range coin will not be fair. This range is called Confidence interval.



Defining some Range is called Confidence Interval

declared by Domain expert.

* We fail to reject the Null hypothesis [within Confidence Interval]

* We reject the Null hypothesis [outside Confidence interval]

S I S T O N S

eg-2 Conclusion

→ Person is Criminal or not of murder case?

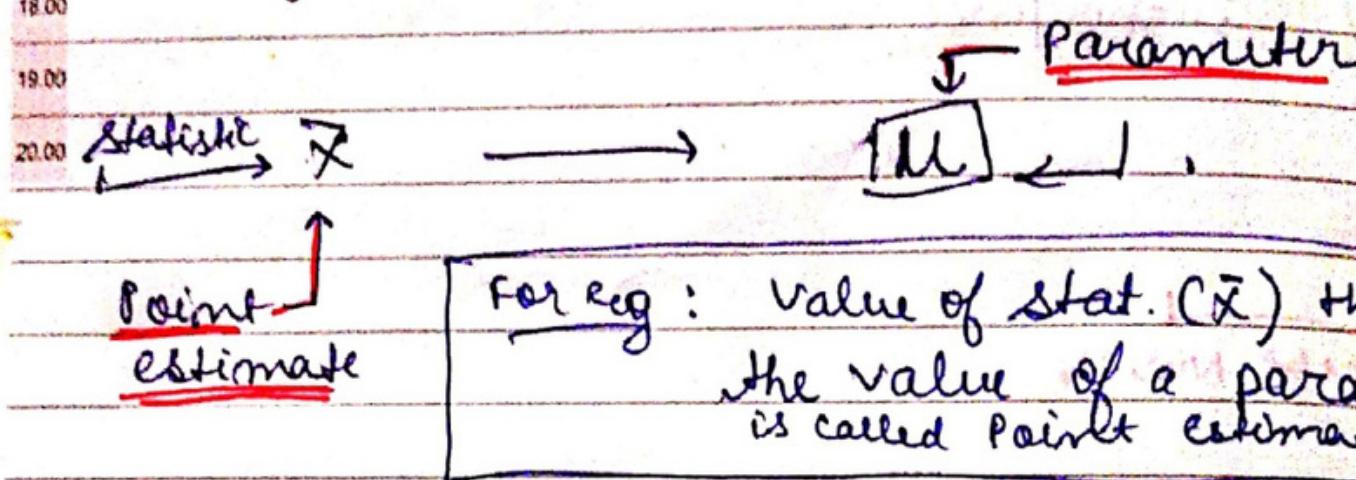
- ① Null Hyp. :- Person is not Criminal
- ② Alternative Hyp. :- Person is Criminal (declared)
- ③ Experiments/Proof :- DNA, finger print, eye-witness
↳ Judge → make decisions.

Monday

44

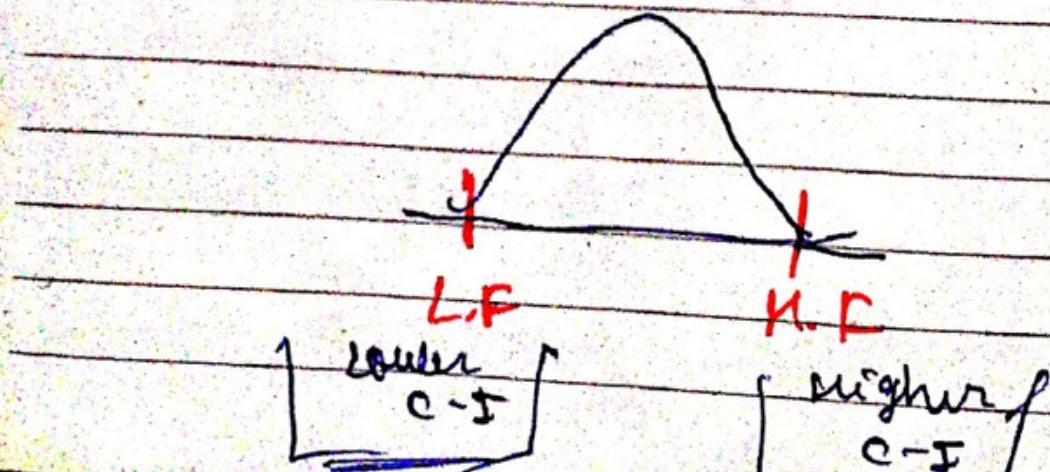
→ Point estimate:- The value of any statistics that estimates the value of a parameter is called Point estimate.

like we are estimating μ (Population mean) on the basis of \bar{x} (sample mean) in inferential stats.



$$\Rightarrow \boxed{\text{Point estimate}} \pm \boxed{\begin{matrix} \text{Margin of} \\ \text{error} \end{matrix}} \rightarrow \boxed{\text{Parameter}} \\ \boxed{\mu \text{ (Population Mean)}} \\ (\geq \text{ or } \leq \text{ than population mean})$$

- Lower fence \equiv Point Estimate, - Margin of error
- Higher fence \equiv Point Estimate, + margin of error



* use always Z -test, when population s.d is given
* use always t-test, when sample s.d is given

■ MAY ■ 2018

142-223/May-21

22

Tuesday

margin of error $\Rightarrow Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

Population s.d. \rightarrow
standard error \rightarrow

25

145-223/May-21

2018 ■ MAY ■

25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100

SCHEDULE

Friday

If in Ques. it is not given if
Sample s.d is given or not
then use formulae

$$\bar{x} \pm t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right)$$

degree of freedom
 $= n-1 \approx 25-1$
 $\Rightarrow 24$

Q6 On the quant test of Amazon exam, a sample of 25 test takers has a mean of 480 with a population s.d. of 85. Construct a 90% of Confidential Interval about the mean.

Sol

- Lower C-I = Point estimate (\bar{x}) - Margin of error

$$\Rightarrow 480 - z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

$$\alpha = 0.1$$

$$\Rightarrow 480 - z_{0.05} \cdot \frac{85}{\sqrt{25}}$$

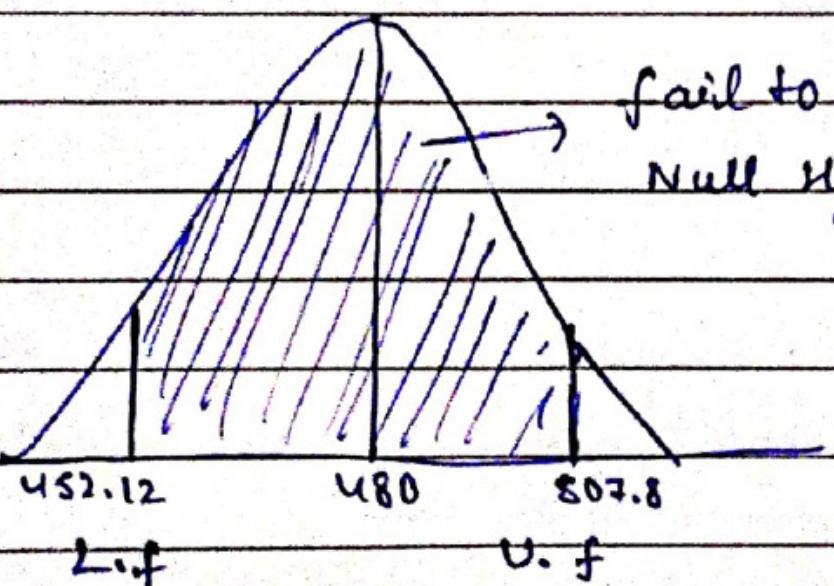
$$\Rightarrow 480 - z_{0.05} \times 17$$

$$\Rightarrow 480 - 1.64 \times 17$$

$$\Rightarrow 452.12$$

- Upper fence $\Rightarrow 480 + 1.64 \times 17$

$$\Rightarrow 507.88$$



fail to reject
Null Hypothesis

Range [452.12 - 507.88]



Monday

09:00

Z-test(if $n \geq 30$, OR population std is given)

10:00

t-test(if $n \leq 30$, s is given)

11:00

12:00

13:00

14:00

15:00

16:00

17:00

18:00

19:00

20:00

→ Hypotheses Testing Prob.

~~Ques~~ A factory has a machine that ~~should~~ ~~must~~ fill 80 ml of Baby medicines in a bottle. An employee believes that avg. amount of baby medicine is not 80 ml. Using 40 sample he measures the avg. amount dispersed by a machine to be 78 ml with a S.D. of 2.5

(a) State Null Hypothesis

(b) At 95% c.I., is there enough evidence to support machine is working properly or not.

Sol:- Null hypothesis, $H_0: \mu = 80$ (machine working properly)

Step I

Alternate H_a, $\mu \neq 80$ (machine not working properly)

Step II

$$\mu = 80 \text{ ml}$$

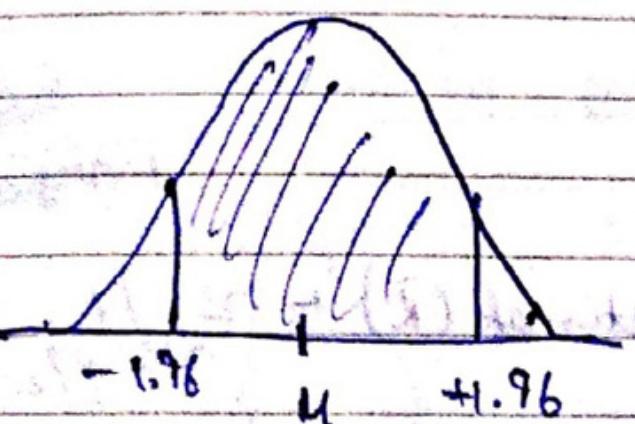
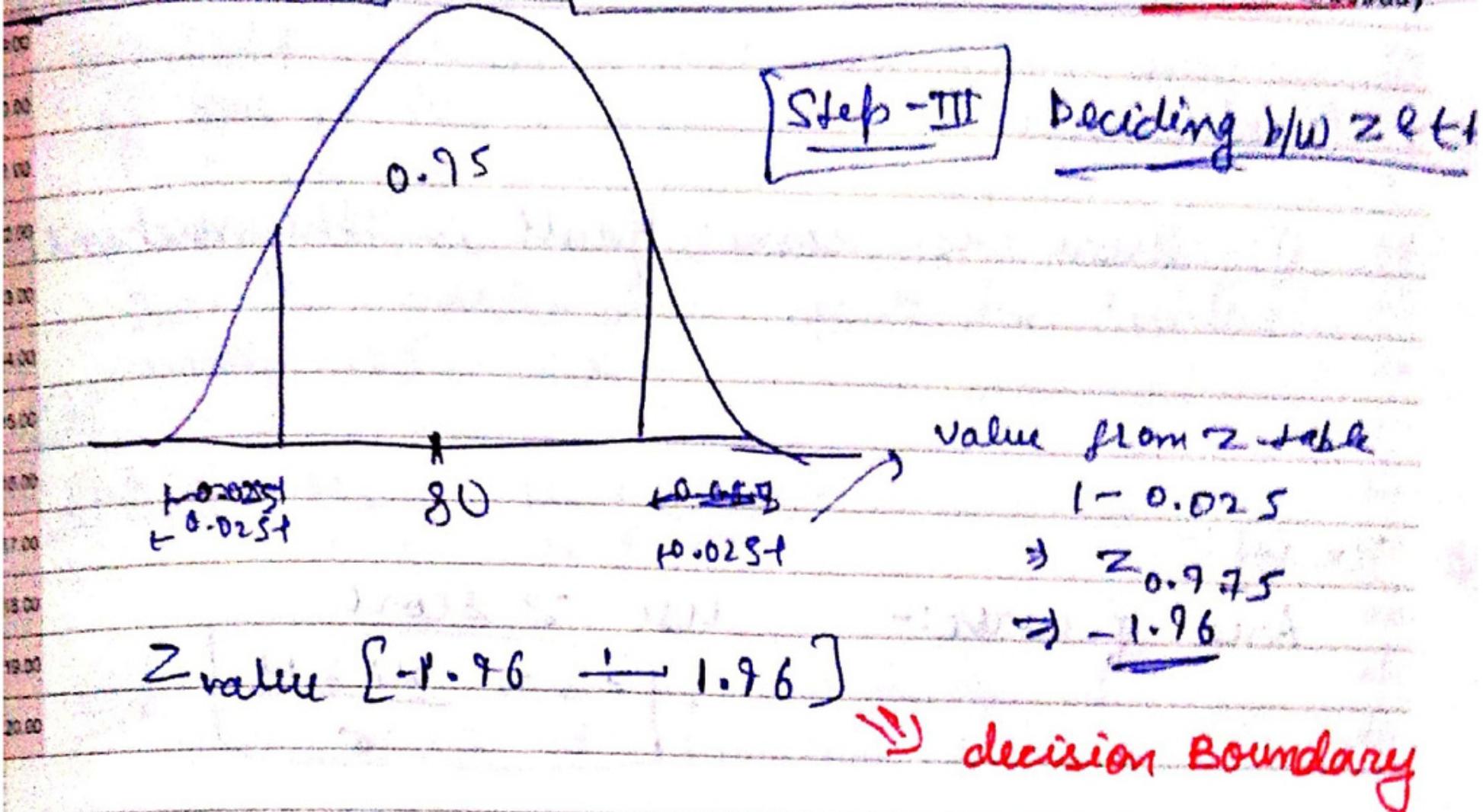
$$n = 40$$

$$\bar{x} = 78 \text{ ml}$$

$$S = 2.5$$

$$\rightarrow C.I. = 0.95$$

$$\rightarrow \alpha, (\beta.v) = 1 - 0.95 \rightarrow 0.05$$



Step IV

→ Now calculate / test statistics (Z-test)

$$z = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \rightarrow 78 - 80$$

$$\frac{2.5}{\sqrt{40}} \rightarrow \text{standard error}$$

$$\Rightarrow -5.05$$

→ Conclusion

My Decision Rule:- If $z = -5.05$ is less than -1.96 or greater than $+1.96$, then we Reject Null H_0 with 95% of Confidence Interval.

⇒ There is some fault in the machine.

2

Q9 A factory manufactures cars with a warranty of 5 years or more on the engine & transmission. An engineer believes that the engine or transmission will malfunction in less than 5 years. He test a sample of 40 cars and finds the average time to be 4.8 years with a standard deviation of 0.50. State the H_0 & H_1 at a 2% of significance level, is there enough evidence to support the idea that the warranty should be revised?

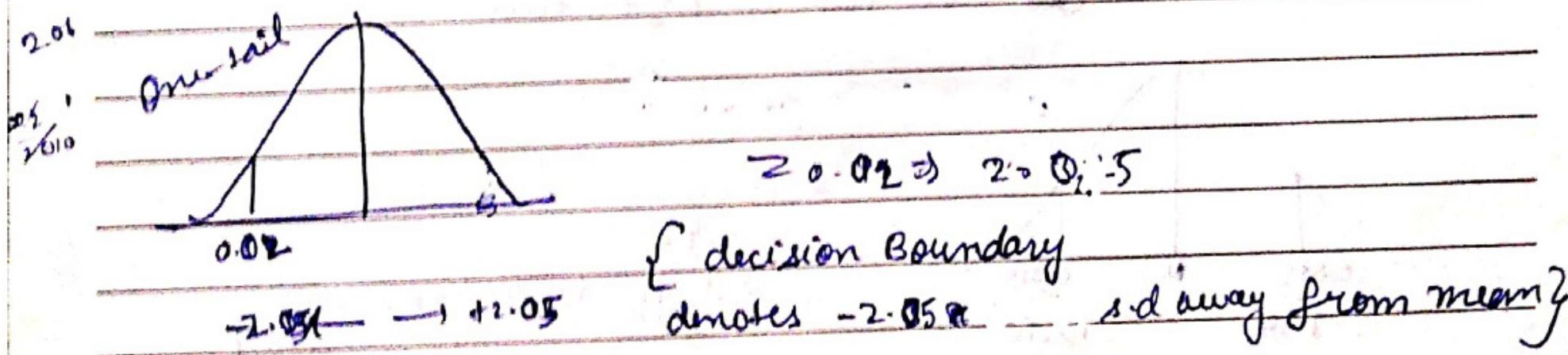
Step 1

$$H_0: \mu \geq 5$$

$$H_1: \mu \leq 5$$

Step 2: $\alpha = 0.02$, $CJ \Rightarrow 0.98$
 $\sigma = 0.5$ (population s.d.)

Step 3: perform ~~Z-test~~ (one-tail) ~~left-tail~~

Step 4

$$Z_{\text{score}} \Rightarrow \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \Rightarrow \frac{4.8 - 5}{0.5/\sqrt{40}} \Rightarrow \frac{-0.2}{0.079} \Rightarrow -2.53$$

Step 4: Conclusion

$$-2.53 < -2.05$$

\Rightarrow Reject null Hypothesis, H_0

\Rightarrow Engine or transmission malfunction is less than 5 years.

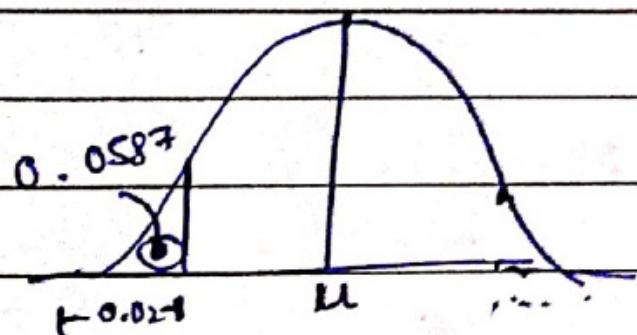
"With the help of z-score whatever area under the curve we are getting is p-value"

Q9 Q9 with the help of P-value

Sol: Step 1: $H_0: \mu = 5$
 $H_1: \mu \neq 5$

Step 2: $\sigma = 0.5$ {one-tail test}
 $\alpha = 0.02$

Step 3: Apply Z-test & find p-value



$$Z_{\text{Score}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \Rightarrow \frac{4.8 - 5}{0.5/\sqrt{40}} = -2.53$$

Area under the curve, $Z_{-2.53} = 0.0587$
so our p-value is [0.0587]

Step 4: $0.0587 (P) < \text{S.V. } (\alpha)$ {Conclusion}

\Rightarrow We reject H_0

\Rightarrow engine malfunction in less than five years.

Q10 A company manufactures bike's batteries with an average life span of 2 years or more years. An engineer believes this value to be less.

Using 10 samples, he measures the average life span to be 1.8 years, with a standard deviation of 0.15.

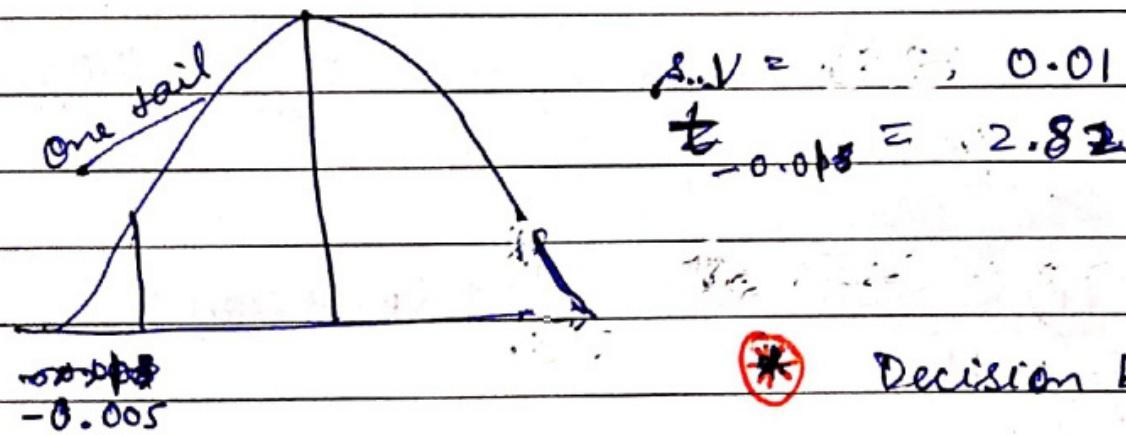
a) State the Null and Alternate hypothesis.

b) At a 99% CI, is there enough evidence to discard the H_0 ?

Sol Step 1: $H_0: \mu_0 \geq 2$
 $H_1: \mu_0 < 2$

Step 2: sample var, $s^2 = 0.15$ {one-tail, left tail
 $\bar{x} = 1.8, n = 10$ (coz eng. believes)

Step 3:- * perform t -test test statistics :-



t -test = ~~df~~ degree of freedom = $n-1 \Rightarrow 9$

$$t_{\text{score}} = \frac{\bar{x} - \mu}{s/\sqrt{n}} \Rightarrow \frac{1.8 - 2}{0.15/\sqrt{10}} \Rightarrow -4.255$$

Step 4: Conclusion



Here (Z score) $-4.255 < -2.82$ (decision Board.)

\Rightarrow Reject H_0

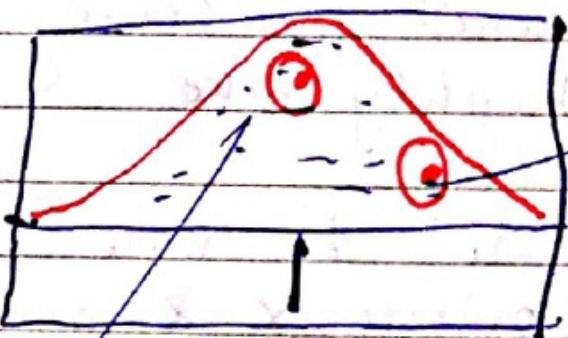
\Rightarrow life span of Battery is less than 2 years **Building Careers Accurately!**

Tuesday

SCHEDULE

09.00 → P-value :- P value \neq Significance Value (?)
10.00
11.00
12.00 chick derived from c-I
13.00
14.00
15.00
16.00
17.00
18.00
19.00
20.00

For eg:- In case of our Mousepad



Here the value $p = 0.02$
that indicates out
of all the hundred
touches we do
randomly 2 touches
will go over here.

High p-value,

let say 0.80

so we can say
out of all hundred touches
we do 80 touches.

* → P value ^{use to} find Hypothesis.

Step 1:- first find Z score. By using
Z score find area under curve using
Z-table.

Step 2:- compare the p-value with S.V.
(the coming value from Z-table is p-value)

If, $p\text{ value} < \alpha\text{ value}$ | $p\text{ value} > \alpha\text{ value}$

∴ Reject null hypothesis | ∴ Accept H_0

Note * if test is 1 tail value is area under curve $\frac{0.05}{2}$
if test is 2 tail value is $0.05 + 0.05 = 0.1$

→ (χ^2) Chi Square Test. [always applied on categorical data]

* Chi Square test claims about population proportions.

* It is a non parametric Test that is performed on categorical data.

* find degree of freedom. ordinal data Nominal data.

* If $\chi^2 >$ Decision Boundary
 ↳ Reject H_0

* there is no 1 tail or 2 tail.

Q1 In a 2000 U.S Census the age of individuals in a small town found to be the following.

<18	$18-35$	>35
20%	30%	- 50%

In 2010, ages of $n = 500$ individuals were sampled. Below are the results.

(out of) 500	<18	$18-35$	>35
	121	288	91

Using $\alpha = 0.05$, would you conclude the population's distribution of ages

has changed in the last 10 years?

Sel:	≤ 18	$18-35$	> 35
Out of 500	observed	expected	expected
	121	288	91
	100	150	250

Step 1: H_0 : The data meets the expected distrib.
 H_1 : The data does not meet " "

Step 2: $\alpha = 0.05$, $CJ = 7.81$.

Step 3: D.O.F {Categories}
 $D.F = \text{Category} - 1$
 $\Rightarrow 3 - 1 \Rightarrow 2$

Step 4: Decision Boundary

* (check in the square table, compare it with df & α value)

: [5.991]

Step 5: Chi Square Test Statistics

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} \text{ OR } \sum \frac{(f_o - f_e)^2}{f_e}$$

$$\Rightarrow \frac{(121-100)^2}{100} + \frac{(288-150)^2}{150} + \frac{(91-250)^2}{250}$$

[$\chi^2 \Rightarrow 232.494$]