

Q1

Q1MD-2

Percapita income of employed women in Jamshedpur

$$PI = \frac{\text{Total Income}}{\text{Tot. No. of employed women}}$$

1) Household survey (census)

we are applying simple mathematical tool ratio. But before doing ratio we need to do lot of work to get data.

So, to get a simple value like this we can't afford big tool based on resources

Statistic is the ^{measure} ~~major~~ of a sample.
Parameter " " " " population.

Descriptive statistics
Inferential statistics $\left\{ \begin{array}{l} \text{statistical estimate} \\ \text{Testing Hypothesis} \end{array} \right.$

Sample \rightarrow A small representative unit of population.

μ = pop. mean

\bar{x} = sample mean

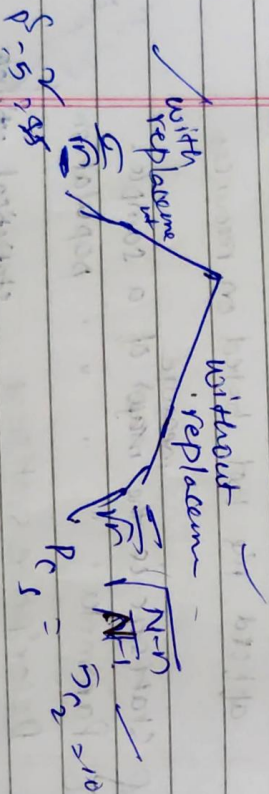
$\sigma \rightarrow$ POP. STD deviation

$s \rightarrow$ sample " " "

$$\sigma^2 = \frac{1}{N} \sum (x_i - \mu)^2 \Rightarrow \sigma = \sqrt{\frac{1}{N} \sum (x_i - \mu)^2}$$

$$s = \sqrt{\frac{1}{(n-1)} \sum (x_i - \bar{x})^2}$$

Standard Error $E(\bar{x}) = \frac{\sigma}{\sqrt{N}}$



Standard error is standard deviation also.

Ex: POP(N) = 2000 $\mu = 66.8$, $\sigma^2 = 2$

10 samples of size = 25

a) with replacement.

$$E(\bar{x}) = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}$$

Mean will be same with and without replacement.

Std. deviation of sampling with replacement

$$= \frac{\sigma}{\sqrt{n}} = \frac{2}{5} > 0.6$$

without replacement.

$$= \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} = \frac{2}{5} \sqrt{\frac{2000-25}{2000-1}} = 0.6$$

$$\frac{D}{N} < 0.05$$

then we can tell size of sample is small

compared to population.

Standard error, $Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$

$$P(66.8 \leq \bar{x} \leq 68.3) = ? \times 100 =$$

$$\bar{x}_1 = 66.8, z_1 = \frac{66.8 - 66.8}{0.6} = -2$$

$\mu_i - x$

Q10/11/12

$\bar{x}_n = 68.3$; $z_{\alpha/2} = 68.3 - 68$
 $\frac{0.6}{0.5}$ > 0.5

$P(-2 \leq z \leq 0.5) = ?$ $0.6867 \times 60 \approx 53$
sample

$P(\bar{x} \leq 66.4)$; $z = \frac{66.4 - 68}{0.6} = -2.66$

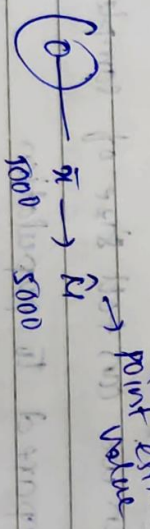
$P(z \leq -2.66) = P(-2.66 \leq z \leq 2.66)$ $\times 80 = 0.0038$
 $\times 80 = 0.304$

Q9

100 stds \rightarrow sample

pop \rightarrow 1546, mean $\mu = 67.45$, $\sigma = 2.93$.

find 95% confidence interval.



mean, $AI = \left[\bar{x} \pm z * SE_{\bar{x}} \right]$

confidence \rightarrow 95%
95% sure

we can create 100 such intervals. In $\frac{1}{100}$ such intervals population mean will be there.

Only 95 intervals are considered as valid.

Just like mean we can have any parameter we can estimate based on the sample statistic.

eg 1

Random sample of 50 BT graders out of total 200 stds mean = 75, $\sigma = 10$. What are 95% confidence limits for estimates of the

mean of 200 graders.

sol 1

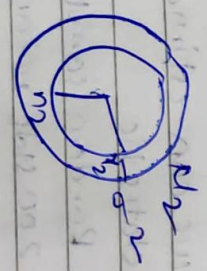
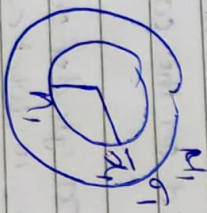
$CI = 75 \pm$ $n(S) = 50$, $n(P) = 200$

$SE = \frac{\sigma}{\sqrt{n}} = \frac{10}{\sqrt{50}} = \frac{10}{\sqrt{200-1}} = \frac{10}{\sqrt{199}} \approx 0.71$

Estimated interval = $\bar{x} \pm Z_{\alpha} \times SE_x$
 $= [72.6 \pm 1.96 \times 1.2227]$

$= [72.6, 74.4]$

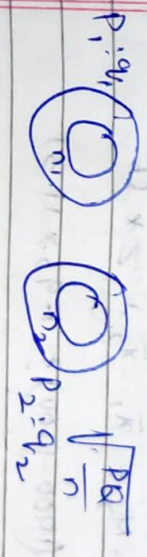
This interval is based on your sample. This is not of such a valid interval. Based on this interval we can estimate population.



$SE = \sqrt{\frac{\sigma_1^2 + \sigma_2^2}{n_1 + n_2}}$
 $\bar{x}_1 = \bar{x}_2 = 1$
 $N_1 = N_2$

$\frac{\sigma_1^2 + \sigma_2^2}{n_1 + n_2}$

eg 1



$P = n_1 P_1 + n_2 P_2$
 $Q = 1 - P$

$SE = \sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$

A sample of 150 brand A bulbs mean = 1400 hrs. $\sigma = 120$ hrs. A sample of 200 brand B lights mean = 1200 hrs, $\sigma = 80$ hrs. Find 95% confidence limits for the diff of mean life time of the population of brands A & B

$SE = \sqrt{\frac{\sigma_1^2 + \sigma_2^2}{n_1 + n_2}} = \sqrt{\frac{120^2 + 80^2}{150 + 200}}$
 $= \sqrt{\frac{6000}{350}} = 11.81$

$$N_1 \sim N_2 = \bar{x}_1 - \bar{x}_2 + Z * SE$$

$$= (1900 - 1500) \pm 1.96 * 11.3$$

$$= 2000 \pm 22.1676$$

$$= [1777.8324, 2222.1676]$$

⑧ random sample of 400 adults and 600 teenagers, 100 adults, 300 teenagers liked it. Construct a 95% confidence limits for the diff in proportion of all adults and teenagers who wa

$$P = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2} = \frac{100 * 100 + 300 * 300}{400 + 600} = \frac{190000}{1000} = 1900 + 1500 = 2220.$$

$$Q = 1 - P =$$

$$SE = \sqrt{220 \left(\frac{1}{600} + \frac{1}{400} \right)}$$

$$P_1 = 0.25 \left(\frac{100}{400} \right) \quad P_2 = 0.5 \left(\frac{300}{600} \right)$$

$$n_1 = 400$$

$$n_2 = 600$$

$$P = n_1 P_1$$

$$= \frac{400 * 0.25 + 600 * 0.5}{1000}$$

$$= \frac{100 + 300}{1000} = 0.4$$

$$Q = 1 - P = 0.6$$

$$SE = \sqrt{(0.4)(0.6) \left(\frac{1}{600} + \frac{1}{400} \right)}$$

$$= \sqrt{0.24 * \left(\frac{1000}{240000} \right)} = 0.031$$

$$P \pm SE = [0.25 + 1.96 * 0.031]$$

$$= [0.19, 0.31]$$

12) axial length of pistons produced is normally distributed with variance of 64mm. Determined sample size such that mean length of pistons can be obtained within a variation of 2mm with a confidence level of 95%.

sol) Variance = 64 \Rightarrow SD $\sigma = 8$ n??

$$P(\bar{x} \pm Z_{\alpha/2} SE_x)$$

$$N - \bar{x} \pm Z_{\alpha/2} \times SE \approx 2 \text{ mm}$$

$$\frac{1.96 \times \sigma}{\sqrt{n}} \approx 2$$

$$\frac{1.96 \times 8}{\sqrt{n}} \approx 2 \Rightarrow n \approx 61.54$$

\rightarrow If the sample mean of 10 elements is 10.

$$x_1 + x_2 + \dots + x_{10} = 10$$

Out of these 10 elements can change

independently so that the mean still remains at 10

Max 9 elements randomly can be changed, both will be depends on sum of 9 elements and SD.

No. of independent components = n-1

No. of " components is called as

" Degrees of Freedom "

$$E(\bar{x}) = \mu = \frac{\sum x_i}{n} = \mu$$

$$\sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

$$E(S^2) = \sigma^2 \quad \sigma^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

If this is replaced by n-1 then

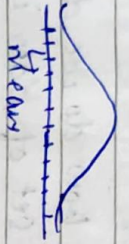
$$E(S^2) = \sigma^2$$

Bessel's correction \rightarrow n is changed to n-1.

P = {2, 3, 4, 5, 6}

N(S) = 2 \rightarrow 2, 2.5, 3, 3.5, 4, 4.5

possible = 1, 2, 3, 4



$$E(x) = \mu$$

Sampling Error with replacement

" " without

No. of samples drawn:

with replacement = $5^2 = 25$ $\sqrt{\frac{\sigma^2}{n}} = SE$

without " " $5^2 = 25$ $\sqrt{\frac{\sigma^2}{n-1}} = SE$

Point Estimation

Ex- POP = 3000, $\mu = 66$, $\sigma = 3$

Sample size = 25 $\bar{x} = ?$, $n = ?$, no. of samples = 80

with replacement = 68, without = 66

$$\sigma = \frac{3}{\sqrt{25}} = 0.6$$

$$\frac{3}{\sqrt{25}} = 0.6$$

This happens b/c size of sample is less when compared to size of population.

If size of the sample by itself $\frac{N}{n} < 0.05$ then we can say size of the sample is small compared to size of population. Then with replacement and without replacement are same.

→ In how many samples of the above problem would you expect to find the mean

a) b/w 66.8 and 68.3 in dev.

Total 80 samples, 68

$$Z = \frac{\bar{x} - \mu}{SE}$$

$$P(66.8 \leq \bar{x} \leq 68.3) = ? \times 80 = ?$$

$$\bar{x}_1 = 66.8, Z_1 = \frac{66.8 - 66}{0.6} = -2$$

$$\bar{x}_2 = 68.3, Z_2 = \frac{68.3 - 66}{0.6} = 0.5$$

$$-2 \leq Z \leq 0.5 \Rightarrow P(-2 \leq Z \leq 0.5) = 0.6687 \times 80 = 53$$

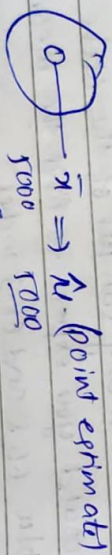
$$P(66.8 \leq \bar{x} \leq 68.3) = 0.6687$$

$$P(-2 \leq Z \leq 0.5) = 0.6687 \times 80 = 53.5$$

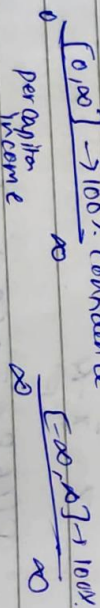
0.30420

② Height of 100 male std 1
 POP = 1546, $\mu = 69.45$, $\sigma = 2.93$

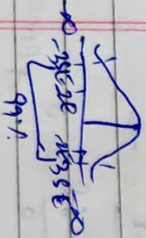
Find 95% confidence interval for estimating the mean height of the std 1.
 point estimate



sample n representative of population.



This interval is meaningless



$$N \pm Z = \bar{x} \pm \frac{\sigma \sqrt{N}}{\sqrt{SE}}$$

$$N \left[\bar{x} \pm Z \cdot SE \right]$$

95%
99% 90%

Based on the confidence level we can find the interval

If the confidence is 99% then 1% is incorrect

$$\hat{\mu} = \left[\bar{x} \pm Z \cdot SE \right] \quad \sigma = 2.93$$

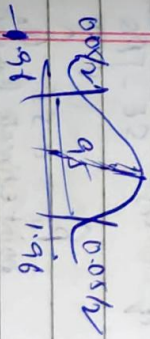
$$= [69.45 \pm 0.98 \times 2.93] \quad Z = \frac{\bar{x} - \mu}{\sigma} = 2.93$$

$$= Z(0.975)$$

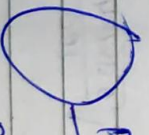
$$SE = \frac{2.93}{\sqrt{100}} = \frac{2.93}{10} = 0.293$$

$$N = [69.45 \pm 0.2834 \times 1.96] \rightarrow \text{confidence of 95\%}$$

$$Z = 1.96$$

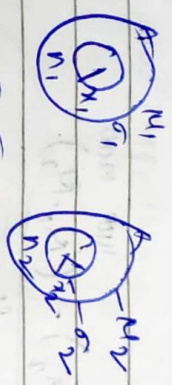


proportional

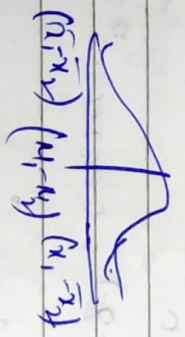


200 smokers
 1000 population
 $P(S) = \frac{200}{1000} = 0.2$
 $\hat{p} = 1 - P(S) = 0.8$
 PR & represent proportions

2 populations



\bar{x}_1, \bar{x}_2 is a new statistic and they also follow a distribution



$$E(\bar{x}_1 - \bar{x}_2) = \mu_1 - \mu_2$$

$$SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

size of sample of pop, and pop

pop

Manufacturer \rightarrow mean = 1400, $\sigma = 200$.

B \rightarrow " = 1200, $\sigma = 100$

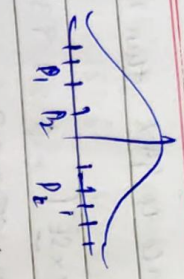
Random sample of 125 (A) 125, (B)

what is prob that brand A bulbs will have a mean life time that is at least.

a) 160 hrs more than the brand B bulbs.

$$SE = \sqrt{\frac{(200)^2}{125} + \frac{(100)^2}{125}} = \sqrt{\frac{40000}{125} + \frac{10000}{125}} = \sqrt{\frac{50000}{125}} = \sqrt{400} = 20$$

Pb3 Sample of 100 voters chosen, 55% of them are in favor of particular candidate.



$$SE_p = E(p_i) = p$$

$$SE_p = \sqrt{\frac{pq}{n}}$$

$$SE_p = \sqrt{\frac{(0.55)(0.45)}{100}}$$

all voters

$n = 100$ $p = 0.55$ $q = 0.45$ $SE = \sqrt{\frac{pq}{n}}$

(estimate) $\hat{p} = p \pm z_{\alpha/2} SE$

$$= 0.55 \pm SE$$

$\hat{p}_0 = 0.55$ $\hat{p}_0 = 0.45$ point estimate for pop. proportion

$$\left[0.55 \pm 1.96 \times 0.0497 \right]$$

$$SE_p = \sqrt{\frac{(0.55)(0.45)}{100}} = 0.0497$$

$$Z = \frac{\bar{x} - \mu}{\frac{SE}{\sqrt{n}}} = \frac{100 - 100}{\frac{SE}{\sqrt{100}}}$$

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{SE_{\bar{x}_1 - \bar{x}_2}}$$

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - 200}{\frac{SE}{20}} = \frac{160 - 200 - 200}{\frac{SE}{20}}$$

$$P(Z) = P(-2) = 0.0227$$

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{SE_{\bar{x}_1 - \bar{x}_2}} = \frac{7 - 8 - 1.38}{\frac{SE}{1.38}} = \frac{-1.38}{0.22} = -6.27$$

$$SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = 1.38$$

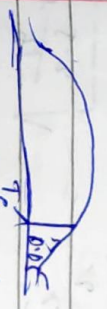
Sample doesn't give evidence to reject null hypothesis

Mean weight = 100gms sample of 25

$\mu = 99$, variance = 4 claim of 0.05

Null hypothesis $H_0: \mu = 100$

$\mu = 100$



Variance $\sigma^2 = 4$
sample size $n = 20$

So, we have to take T-distribution

Degree of freedom = $n - 1 = 20 - 1 = 19$

critical $t = 1.7311$

$$t = \frac{\bar{x} - \mu}{\frac{SE}{\sqrt{n}}} = \frac{97 - 100}{\frac{1.6}{\sqrt{20}}} = \frac{-3}{0.3577} = -8.41$$

So, it doesn't fall in rejection.

$T < T_c$ failed to reject. The claim is false.

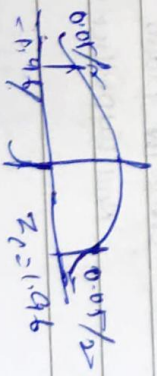
diameter (bit) = $b_2 \rightarrow$ sample mean = 40

sample $X_6 \rightarrow$ Mean = 38.5
Var = 2.5

var = 4

$H_0: \mu_1 = \mu_2$

$H_1: \mu_1 \neq \mu_2$



$$t = \frac{\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)}{SE} \quad | \quad SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$SE = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} = \sqrt{\frac{2.5}{\frac{16}{80} + \frac{1}{20}}}$$

$$SE = SP \sqrt{\frac{1}{16} + \frac{1}{20}} = 0.59$$

$$\frac{2.5 - 4}{0.59} = \frac{1.5}{0.59} = 2.542$$

If falls in the rejection area of null hypothesis, reject the null hypothesis. Managers claim is wrong.

→ 90% sample = 80, it cured 68% patients. check the claim at 0.05. Null hypothesis - $H_0: P \leq 90\%$. population proportion

All hypothesis $H_1: P > 90\%$.
 Standard error.
 $SE = \sqrt{\frac{pq}{n}} = \sqrt{\frac{0.9 \times 0.1}{80}}$

$$Z = \frac{P - p}{SE} = \frac{0.68 - 0.9}{\frac{0.034}{80}} = -0.47$$

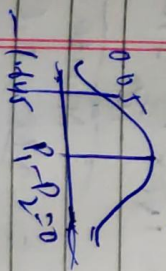
Not in the rejection area then failed to reject null hypothesis.

→ $R_2 > R_1$,

sample 1 = 110 days | 52 = 130 days
on 90 days target > | 112 days target >

check at 0.05 level

Null $H_0: P_2 \leq P_1$, $H_1: P_2 > P_1$



$$SE = \sqrt{\frac{p_1 q_1 + p_2 q_2}{n_1 + n_2}} = \sqrt{\frac{0.9 \times 0.1}{110} + \frac{0.9 \times 0.1}{130}}$$

→ $P_2 > 1 - P_1$
 $P_2 > P_1 + \frac{1}{2}$

$$SE \approx 0.04 \quad Z = \frac{(P_1 - P_2) - (P_1 - P_2)}{SE}$$

$$= \frac{0 - \left(\frac{110}{110} - \frac{112}{130} \right)}{0.04} \approx -0.922$$

95% not in rejection area.

So, we can't reject null hypothesis.

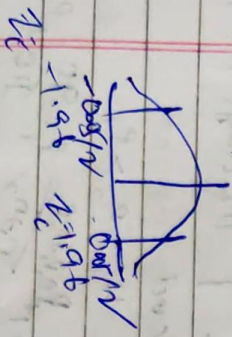
→ Sample 1

City X → 125 out of 150 sig. level = 0.05
 City Y → 133 out of 180.

Let $H_0: P_1 = P_2$
 $H_a: P_1 \neq P_2$

proportion of EG holders from X, Y

Sampling distribution of proportion



$$Z = \frac{\hat{p} - p}{SE} = \frac{p_1 - p_2}{\sqrt{\frac{p_1 + n_1 p_1}{n_1 + n_2} + \frac{p_2 + n_2 p_2}{n_2 + n_2}}}$$

$$= \frac{150 \times \frac{125}{150} + 180 \times \frac{133}{180}}{150 + 180}$$

$$= \frac{268}{330} \approx 0.81$$

$$\alpha = 0.14$$

$$SE = \sqrt{p \hat{p} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \approx \sqrt{0.88 \times 0.14 \left(\frac{1}{150} + \frac{1}{180} \right)}$$

$$\approx 0.04$$

$$Z = \frac{P_1 - P_2}{SE} = \frac{\left(\frac{125}{150} - \frac{133}{180} \right)}{0.04} = 1.$$

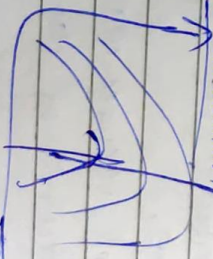
$$= 1.99$$

z lies within range.

Null hypothesis is not rejected

Failed to reject null hypothesis

t-distribution



pop. when standard dev is not given and pop. size ≤ 30 .
 If $n > 30$ then it is similar to z-dist.

highly skewed range $[0, \infty]$

Chi-square and F-test

F-test we use for 2 populations.

Calc: $F = \frac{S_1}{S_2}$
Compare variance of 2 POP by taking 2 samples

Chi-square is also testing variance w.r.t. one population.

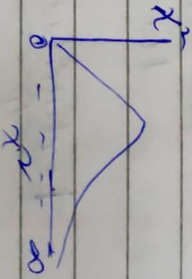
In F-test of testing inferential hypothesis either with 1 population or 2 population

Row \rightarrow S.D of sample

level \rightarrow denom of $F = \frac{S^2}{\sigma^2}$

well temp with row F -test start from 0 to ∞ .

highly skewed distribution



we will tell sample variance w.r.t. pop. variance.

Other way of representation of chi-square.

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

observed frequency

expected frequency

Expected frequency

This way of chi-square is used for goodness of fit.

No. of car accidents per month in a certain city were as follows:
12, 18, 20, 2, 14, 10, 15, 6, 9, 4

Are these frequencies in agreement with the belief that accident conditions were same during this 10 month period.

Null hypothesis H_0 : Accident condition were same

H_1 : conditions were not same.

total No. of accidents = 110.

If condition same no. of accidents/month = 11

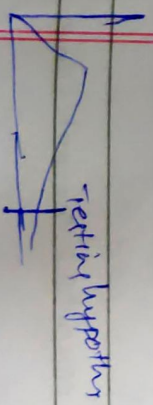
Based on this

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

$O_i - E_i > 0$ & < 0

if the diff is significantly

high then χ^2 value will be high



If χ^2 is beyond certain limit we can tell the alt. hypothesis is true & if it is within the limit, we can tell that we can't reject the null hypothesis.

$$\chi^2 = (0-11)^2 + (18-11)^2 + (20-11)^2 + (2-11)^2 + (14-11)^2 + (10-11)^2 + (15-11)^2 + (6-11)^2 + (9-11)^2 + (4-11)^2$$

$$= \frac{1+49+81+81+9+1+16+25+4+49}{11}$$

$$= \frac{2 \times 81 + 2 \times 49 + 16}{11} = \frac{316}{11} = 28.7$$

Degree of freedom = $n-1 \rightarrow 9$ 10 observations

At 5% level of significance of 1 D.F. $\rightarrow 9$ then

$$\chi^2 = 16.19$$

We can't reject null hypothesis

$$(28.7 > 16.19)$$

In the accident conditions are not same.

Based on the data we are checking whether our assumption is correct or not. This is good fit.

Q22 Theory predicts proportion of items in 4 groups A, B, C, D \rightarrow 9:3:3:1. In an exp among 1600 items, the no's in grp were 882, 313, 287 and 118. Does this exp support theory.

Q21 Null hypothesis H_0 :

H_1 :

$$\text{exp'd} = \frac{\text{Total}}{4} = \frac{1600}{4} = 400$$

$$\chi^2 = (882-400)^2 + (313-400)^2 + (287-400)^2 + (118-400)^2$$

$$= \frac{2529324 + 7569 + 12769 + 79524}{400}$$

$$= 8304.65$$

sol 2 Null hypothesis H_0 : Experiment supports the theory

H_1 : Exp doesn't support the theory.

Theory $\Rightarrow A:B:C:D = 9:3:3:1$

Experiment T_1 expected (9:3:3:1)

A	82	900
B	313	300
C	282	300
D	118	100
		<u>1600</u>

$$\chi^2 = \frac{(82-900)^2}{900} + \frac{(313-300)^2}{300} + \frac{(282-300)^2}{300} + \frac{(118-100)^2}{100}$$

$c = 4, d = 3$



DOF = 3
At 5% significance level, critical point = 7.815

$4.9 < 7.815$

int: Can't reject null hypothesis

The experiment support theory.

PB 3-1

Records taken of no of male and female births in 800 families having 4 children are given as:

male	F	families
0	4	32
1	3	128
2	2	290
3	1	236
4	0	64

Test whether the data is consistent with the hypothesis that male & female births are equally likely.

sol 3

Null H_0 : M & F births are equally likely

H_1 : Not equally likely

Fam Expected

32	160
128	160
290	160
236	160
64	160

prob. of male births is same as prob. of female births $\therefore H_0$

$$X^2 = 108^2 + 18^2 + 180^2 + \dots$$

$$P(m) = P(t) = 1/2 \quad \text{Total families} = 800 \quad N$$

In a family of 4 children \rightarrow No. of trials = 4
 Thus may be 4 male, 3/2/1/0 male

Binomial distribution
 $P(X=r) = {}^N C_r \cdot p^r \cdot q^{N-r}$

$N=4$
 $P(X=0) = {}^4 C_0 \cdot (1/2)^0 \cdot (1/2)^4 = 0.0625$

$P(X=1) = {}^4 C_1 \cdot (1/2)^1 \cdot (1/2)^3 = 0.375$

$P(X=4) = {}^4 C_4 \cdot (1/2)^4 \cdot (1/2)^0 = 0.0625$

Prob. that there is 1 male child.

Actual Exp. no. of families

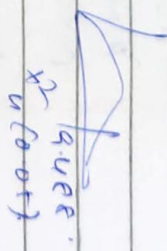
32	800 x 0.0625 = 50
128	800 x 0.25 = 200
240	800 x 0.375 = 300
236	800 x 0.25 = 200
64	800 x 0.0625 = 50

$$X^2 = \frac{18^2}{50} + \frac{22^2}{200} + \frac{10^2}{200} + \frac{38^2}{200} + \frac{14^2}{50}$$

$$= 6.48 + 2.42 + 0.333 + 6.98 + 3.92$$

$$= 19.63$$

Critical Point = 9.488



Reject null hypothesis that male & female births are equally likely.

Prob 51
 The diet of typing mistakes committed by a typist is given as:

Mistake/page	0	1	2/3	4	5
No. of pages	142	156	69/29	5	1

Assuming that diet to be random find out the

exp. no. of pages containing 0, 1, 2, 3, 4, 5 mistakes resp. Test a hypothesis that there is no significant difference among the observed & expected no. of pages containing the mistakes.

Q11

H_0 :
 H_a :

Poisson distribution,

$$P.M.F = \frac{e^{-\lambda} \cdot \lambda^x}{x!} \quad r = 0, 1, 2, \dots, \infty$$

λ = Avg. rate of occurrence.

$$\begin{aligned} \text{Tot. No. of mistakes} &= 5 \times 1 + 4 \times 5 + 3 \times 20 + 2 \times 89 + \\ & \quad 1 \times 156 = \end{aligned}$$

$$= 5 + 20 + 81 + 138 + 156 = 400$$

Tot. Page = 400.

Mistake/page = $1 \rightarrow \lambda$.

$$P(X=0) = \frac{e^{-1} \cdot 1^0}{0!} = 0.3679$$

$$P(X=1) = \frac{e^{-1} \cdot 1^1}{1!} = 0.3679$$

$$\frac{1!}{1!} = 0.184$$

$$= 0.0673, 0.015$$

$$P(X=5) =$$

$$\frac{e^{-1} \cdot 1^5}{5!} = 0.003$$

Expected no. of page containing 0 mistake = 148.

3679

$$\chi^2 = 1.45$$

$$\chi^2_{(5)} = 11.07$$

The no. of mistakes that occur randomly can be calculated using Poisson dist. So, Poisson dist. is best fit.

Q11/12

If there is a significance diff b/w expected level and the given (observed) level.

If observed \geq expected then chi-square will be 0. If the diff is more and more chi-square will be higher.

chi-square one handed test (right handed test)

Test of independence

$H_0: P_1 = P_2$ (based on proportion)

$H_a: P_1 \neq P_2$

ⓐ ⓑ

85 $H_0: P_1 = P_2 = P_3 = P_4$

(Multiple population in terms of proportion)

$H_a: P_1 \neq P_2 \neq P_3 \neq P_4$

Proportions are not same in at least 2 pops
you can test multiple times taking 2 pops at a time.

The prob of null hypothesis is high
on this case we can use chi-square test.

eg:-

	NE	South	Central	West coast	Total
Present Method	68	75	52	39	234
New Method	32	41	33	31	141
Total	100	120	90	110	420

Area is divided into 4 populations. Students applied for interview Agency is thinking about new method of selecting candidate. Sample sizes are 100, 120, 90, 110

In first sample 68 out of 100 choose present method and 32 choose new method
sample proportion = $\frac{68}{100}$ present, $\frac{32}{100}$ new

From samples we can only study proportion.

H_0 : Proportion of students regarding the method is same across regions ($P_1 = P_2 = P_3 = P_4$)

$H_a: P_1 \neq P_2 \neq P_3 \neq P_4$

If we aggregate all the samples then we will get sample for population.

like $100 + 120 + 90 + 110 = 420 \rightarrow$ sample size.

Total no. of candidate referring present method = 229, new method = 191.

Proportion of candidate preferring present method = 0.664 , new method = 0.336

These proportions can be used as an estimate

proportion of candidates.

So, in NE region,

Present: 66.4% \leftarrow expected frequency
 New: 32.6% \leftarrow actual frequency

SE	75	79.7
NS	40.3	

Central.

52	59.8
33	30.2

West Coast	29	73.1
	31	26.9

Chi Square, $\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$

$$= \frac{(66 - 66.4)^2}{66.4} + \frac{(32 - 33.6)^2}{33.6} + \frac{(75 - 79.7)^2}{79.7} + \frac{(40.3 - 32.6)^2}{32.6}$$

$$= \frac{(45 - 40.3)^2}{40.3} + \frac{(59 - 59.8)^2}{59.8} + \frac{(33 - 30.2)^2}{30.2} + \frac{(31 - 26.9)^2}{26.9}$$

$$= 0.0388 + 0.026 + 0.27976 + 0.5481 + 0.1815$$

$$= 0.2591 + 0.426 + 0.945$$

$\chi^2 = 2.26$

As here we have table Degrees of Freedom will be different.

As per the first row $DOF = 4 - 1 = 3$
 " " " " column, $DOF = 2 - 1 = 1$

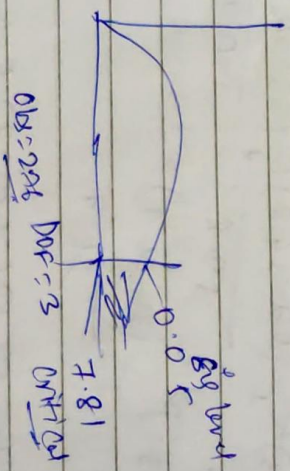
Notable consist of $4 \times 2 = 8$ elements.

DOF for this table = $3 \times 1 = 3$

In this case

$$H_0 = P_1 = P_2 = P_3 = P_4$$

$$H_a = P_1 \neq P_2 \neq P_3 \neq P_4$$



Failed to reject null hypothesis.

Proportion of candidates in their regions are not significantly different. Preference of methods is independent in these regions.

Two attributes, one is method and other is region. Here both the attributes don't influence the preference.

PB 4

hrs spent

Avg. grade.

	A	B	C	D	E	Total
< 5 hrs	13 60.825	10	11	16	5	55
5-10 hrs	20 11.825	22	22	19	2	85
10-15 hrs	9 19.325	22	21	16	32	155
> 20 hrs	8 11.825	11	41	24	11	95
Total	50	75	150	72	50	497

$$H_0 : P_1 = P_2 = P_3 = P_4 = P_5$$

$$H_1 : P_1 \neq P_2 \neq P_3 \neq P_4 \neq P_5$$

Proportion of Ads who listen to music

$$P_1 = \frac{55}{400} = 0.1375 \quad P_2 = 0.2375 \quad P_3 = 0.3875$$

$\xrightarrow{95/400}$
 $\xrightarrow{185/400}$

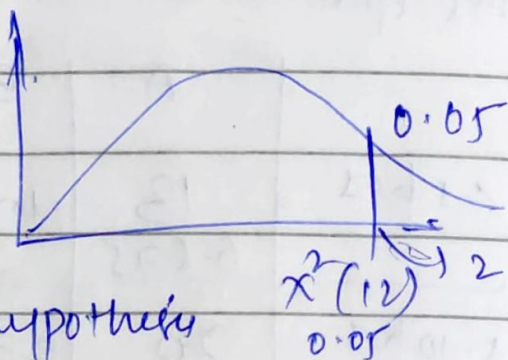
$$P_4 = 0.2375 \rightarrow 95/400$$

$$\chi^2 = 63.83 \quad \text{Dof} = (5-1) \times (4-1) = 12$$

5% sig

Critical $\chi^2 = 21.03$

$\chi^2 > \text{critical}$
 reject the null hypothesis



$$P_1 \neq P_2 \neq P_3 \neq P_4 \neq P_5$$

Alt. hypothesis is true

Proportions are dependent.

Getting avg. grade is dependent on no. of hrs listening to music

One attribute (getting grades) dependant on other attribute (listening to music).

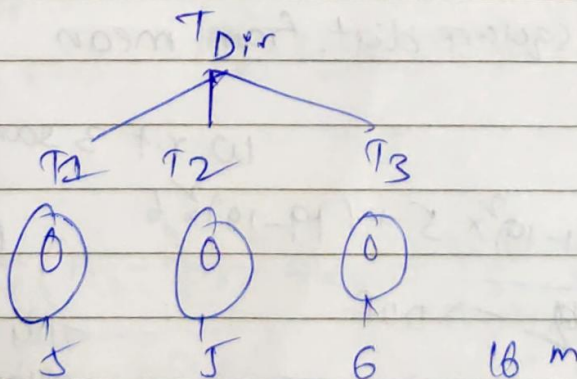
Suppose if we have 2 populations A and B when we can we say both the populations are same? when mean and variance together are same (std. dev)

To test for the significance of the diff among more than 2 sample means.

Training director wanted to evaluate 3 different training methods to determine whether there were any differences in effectiveness of the training methods.

After completion of the training period, she chose 16 new employees assigned at random to 3 training methods.

Counting the production O/P by 16 trainees, she summarized the data and calculated the mean production of the trainees.



To determine the grand mean, of the method can be followed.

Grand Mean = $\frac{1}{5} (15+18+19+22+19) + \frac{1}{5} (22+29+18+21+15)$
 $= 19$

Method 1	$\frac{N_1}{N}$	$\frac{N_2}{N}$	$\frac{N_3}{N}$
	$\frac{1}{5}$	$\frac{2}{5}$	$\frac{1}{5}$
	$\frac{1}{5}$	$\frac{2}{5}$	$\frac{1}{5}$
	$\frac{1}{5}$	$\frac{2}{5}$	$\frac{1}{5}$

Grand Mean = 19

\bar{x}_1 (Mean of Sample 1) = 19 $\bar{x}_2 = 21$ $\bar{x}_3 = 19$

$H_0: \mu_1 = \mu_2 = \mu_3$



Grand mean = 19

Variance \rightarrow Avg. square dist from mean

$$= \frac{\text{Total sq. Dist}}{\text{number}}$$

$$= \frac{5(19-19)^2 + (21-19)^2 \times 5 + (19-19)^2 \times 5}{20}$$

UDF + 3 sample mean

DOF = 3-1 = 2

Total distance of 3 samples we have taken element in terms of samples

$$\frac{20+20}{2} = 20$$

This can be an estimated variance in terms of sample mean

29/10/21

ANOVA

Each of the sample is drawn from a normal population and that each of the population has the same variance.

If the sample size is large - Normality assumption is not required.

If null hypothesis is true, classifying into 3 columns is unnecessary.

3 samples will have their own sample variance,

S_1^2	S_2^2	S_3^2
15, 18, 19, 22, 11	22, 29, 18, 21, 12	18, 29, 19, 16, 22, 15
Mean = $\frac{95}{5} = 19$	Mean = $\frac{107}{5} = 21$	Mean = $\frac{107}{5} = 21$

$$S_1^2 = \frac{2^2 + 1^2 + 2^2 + 25 + 36}{4} = 17.5$$

$$S_2^2 = \frac{1^2 + 8^2 + 9^2 + 10 + 16}{4} = 12$$

$$S^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2 + (n_3-1)S_3^2}{(n_1-1) + (n_2-1) + (n_3-1)}$$

$$= \frac{4 \times 17.5 + 4 \times 12 + 5 \times 12}{14} = 14.714$$

First one, we got $s^2 = 20$.

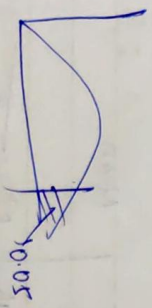
Estimated variance we got based on sample mean.

second one, variance within the samples.

When populations are same then there is no significant difference between their 2 variances (var. among sample mean and var. within among samples).

We can use F-test to compare 2 variances. So, that we can know the difference between variances.

$$F_{\text{test}} = \frac{\sigma_1^2}{\sigma_2^2} = \frac{20}{14.769} \rightarrow 1.354$$



$$F(20, 13) = 3.81 \rightarrow \text{critical}$$

$$1.354 < 3.81$$

Failed to reject null.

Populations are same. \rightarrow This is from F-test.

Variances are same \rightarrow This is assumption.

i.e., means are same $\therefore N_1 = N_2 = N_3$.

Even if we have 20 populations to compare, those 20 populations can be converted into 2 variances i.e., variance amongst sample means and the other is variance within samples.

Problem McDonald,

P1	3	4	4.5	3.5	4	4	Arg
P2	3	3.5	4.5	4	4.5	4.2	
P3	2	3.5	5	6.5	6	4.6	
P4	3	4	5.5	2.5	3	3.6	
							$(\bar{x}) = 4.075$

$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$ $H_1: \mu_1 \neq \mu_2 \neq \mu_3 \neq \mu_4$

Variance among samples.

$$= \sqrt{\frac{4(0.2)^2}{3} + \frac{5(4.075 - 4.1)^2}{3} + \frac{5(4.075 - 4.6)^2}{4} + \frac{5(4.075 - 3.6)^2}{3}}$$

$$= \sqrt{\frac{0.028125}{3} + \frac{0.003125}{3} + \frac{0.000625}{4} + \frac{0.2256}{3}}$$

3

$$s^2 = \frac{2.5875}{3} = 0.8625$$

$$s^2 = 4 \left(\frac{0.000625}{4} \right) + 4 \left(\frac{0.000625}{4} \right) +$$

$$S_1^2 = \frac{(4-3)^2 + (4-4)^2 + (4-5)^2 + (4-3.5)^2 + (4-4)^2}{4} = 0.875$$

$$S_2^2 = 0.925$$

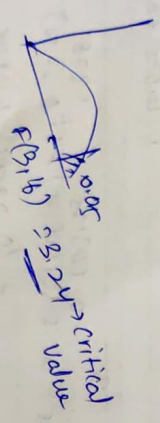
$$S_2^2 = 0.925$$

$$S_2^2 = 1.425$$

$$\sigma_1^2 = \frac{4(0.875) + 4(0.925) + 4(1.425) + 4(1.425)}{16} = 2.21 < 1.6625$$

$$F\text{-value} = \frac{\sigma_1^2}{\sigma_2^2} = \frac{0.875}{1.6625} = 0.526 \approx 0.508$$

DOF = 16



$$F = \sigma_1^2 = \sigma_2^2$$

We can not reject null hypothesis.
There is no sig. diff b/w estimated population variance.
Assumption → variance are same. Do, mean are same.
∴ Avg. service time is same.

Non-Parametric Tests

Also applicable to a situation where data points instead of actual values are as ranks.
It can also apply to assumptions.

eg: SAT scores of students of 2 state universities.

University A →	1000	1100	600	700	1200	950	1050
	1250	1400	850	1150	1200	1500	
	600	700					
B →	920	1120	830	1360	650	720	890
	1600	900	1140	1550	550	1240	940

$H_0: \mu_1 = \mu_2$ (Assumption → variances are same)
 $H_1: \mu_1 \neq \mu_2$

Without assumption we can have parametric test i.e. U test (You test).
Sample elements should be expressed in terms of ranks.

Aggregate both universities and rank all the samples.
 $U_1 \vee A \rightarrow 16$ 18 .

$R \rightarrow$
What is the total rank of all the elements for both universities R_1 and R_2

$n_1 = 15$ $n_2 = 15$

$R_1 = 246$ $R_2 = 218$

Based on these values, $U = \frac{n_1 n_2 + n_1(n_1+1)}{2} - R_1$

$$U = \frac{n_1 n_2}{2} ; SE = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$$

$U = \frac{15 \times 15}{2} + \frac{15(16)}{2} - 246$

$= 225 + 120 - 246 = 99$

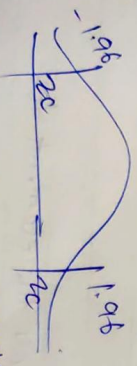
Mean, $M_U = \frac{225}{2} = 112.5$, $SE = \sqrt{\frac{15 \times 15 (31)}{12}} = 24.1$

U distribution can be approximated as Z -distribution when $n_1, n_2 > 10$. Here we need to refer U table

Small samples will be used in biology

$Z = \frac{U - M_U}{SE} \Rightarrow Z = \frac{99 - 112.5}{24.1} = -0.5601$

Failed to reject null hypothesis
-0.601 is b/w critical value



U test (Non-parametric) doesn't take assumptions

We can also write as

$U = \frac{n_1 n_2 + n_2 (n_2 + 1)}{2} - R_2$

$= \frac{225 + 120 - 218}{2} = 129$

$Z = \frac{129 - 112.5}{24.1} = 0.60$

11/12

Kruskal-Wallis (K) statistic
2 assumptions needed regarding

- 1) Normal distribution
 - 2) Variance is same for all the population
- Then we cannot do the ANOVA.
This will do if we have more than 2 pop. and the data is converted to rank data.

Eg:- written examination score of student by 3 diff methods

Video Cassette : 74 88 82 93 55 70 → $n_1 = 6$
 $R_1 = 41$

Audio Cassette : 78 80 65 57 89 → $n_2 = 5$

$R_2 = 3$

Class room : 88 83 50 91 84 77 94 81 92 → $n_3 = 9$

$R_3 = 4$

$n = 6 + 5 + 9 = 20$

K-statistic

$$K = \frac{12}{n(n+1)} \sum_{j=1}^k \frac{R_j^2}{n_j} - 3(n+1)$$

In this example $j=1, 2, 3$.

$$K = \frac{12}{20 \times 21} \left(\frac{(6)^2}{5} + \frac{(6)^2}{6} + \frac{(6)^2}{9} \right) - 3(21)$$

$H_0: N_1 = N_2 = N_3$
 $H_1: N_1 \neq N_2 \neq N_3$

of K -statistic = 11 U.S.

Just like U-statistic we also have K -statistic. If the size of the sample is ≤ 5 then normal distribution is used.

If the size of sample > 5 then chi-square dist. is used. If we need to verify hypothesis

$n_1 + n_2 + n_3 = n \Rightarrow D.O.F = 2$

$\chi^2_{(2, 0.05)} = 5.991$

Failed to reject null hypothesis.

Correlation and Regression

Degree of similarity b/w 2 variable. Each variable can be expressed in terms of data.

They are called as univariate statistic.

Eg:- Mean, Median, Variance, S.D, mode etc.

Multivariate statistic.

We can have multiple variates explained from same data.

Eg:- Correlation, Regression.

There are measures for multivariate.

Degree of similarity b/w variate can be seen.

In simple terms Degree of Similarity (DOS) can be explained in linear

If we explain in terms of linearity of by regressing a straight line then we will call as linear correlation.

If can be explained with Pearson linear coefficient.

Pearson coefficient is covariance of 2 variates

If can be explained in terms of

$$r_c = \frac{\text{Cov}(X, Y)}{\sigma_x \cdot \sigma_y}$$

2 sets of data X, Y

$$\sigma_x^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

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

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

If $R_c = -1$ then

2 variates are perfectly -vely correlated

i.e. If one variate is at max then other is at min.

At $R_c = +1$ then

2 variates are perfectly +vely correlated

i.e. If one variate is at max then other is also at max.

$$S_p = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

~~Spearman~~ Spearman Correlation Coefficient

Non-parametric equivalent of Pearson coefficient

X Y

For Pearson both ^{Data Set} variates can be expressed in terms of

ranks and 'd' is difference in the ranks.

n is no. of ranks.

Regression

If we have 2 variates x and y and ~~when~~ we measure the degree of relation

Regression is a functional relation b/w variables.

If 2 variates x and y, can we build a relationship b/w them. i.e. can we express x in terms of y or y in terms of x.

If we express as $y = f(x)$ then y is dependent and x is independent vari.

Eg:- $y = 5x^2$

i.e. if we have a certain value of 'x' we can calculate y value.

Instead of single independent var we can have multiple independent variables

Eg:- $y = 3x_1 + x_2$ - linear relationship.

It can be expressed as either a linear relationship or non-linear form of relationship.

Based on this principle,

If we have multiple variates $x_1, x_2, x_3, \dots, x_n$ over a period of time they got same time. If one of the variate is an dependent variable, then it is expressed as $x_1 = f(x_2, x_3, \dots)$. Then we have to see what kind of relationship.

x_1	x_2	x_3	x_4	x_5
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-

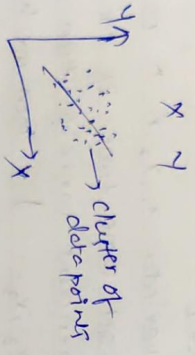
$x_1 = f(x_2, x_3, x_4, x_5)$

If the dependent var is a function of a single variable then the regression is a simple regression.

It deals with 1 independent variable.

If we have multiple no of independent variable then it is known as multiple regression. Likewise, the relationship b/w dependent & independent var than it is a linear regression model. If the relationship is not linear than it is non-linear regression model. All these are estimated based on data points.

Simple Regression:
2 variables
independent variable
Dependent

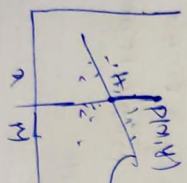


cluster of data points

$$Y = f(x)$$

If we try to explain in linear than we can draw a line then most of the data points are nearby that line passes through most of data point

We need to draw a line such that a point should be on line at the distance b/w the point & line is min.



$$Y = f(x) = a + bx$$

$$x = (y - a) / b$$

$$P_H = [y - (a + bx)]$$

If P is below line

If P depends on the point values and the line

05/11/21

So, we will take square factor of this to make it +ve

$\sum [y - (a + bx)]^2 = E(\text{est})$
E is square distance of the point from the straight line.

It is principle of .
we try to minimize the distance b/w point & line we try to minimize E.

To minimize E the necessary condition is first derivative should be 0. So, E has to be differentiated w.r.t. 2 variables (slope & intercept)

$$\frac{\partial E}{\partial a} = 0$$

$$\frac{\partial E}{\partial b} = 0$$

First order derivatives

$$a + bx - y$$

b → slope → intercept

These 2 first order derivatives will give rise to normal equations in terms of data points

Regression coefficients.

$$b = \frac{\sum xy - n\bar{x}\bar{y}}{\sum x^2 - n\bar{x}^2}$$

$$a = \bar{y} - b\bar{x}$$

a and b are regression coefficients

$$\sum y = na + b\sum x$$

$$\sum xy = a\sum x + b\sum x^2$$

When these equations are solved then we can get regression eqn ⇒ y = a + bx

get regression eqn ⇒ y = a + bx