

Regression Analysis.

classmate

Date _____

Page _____

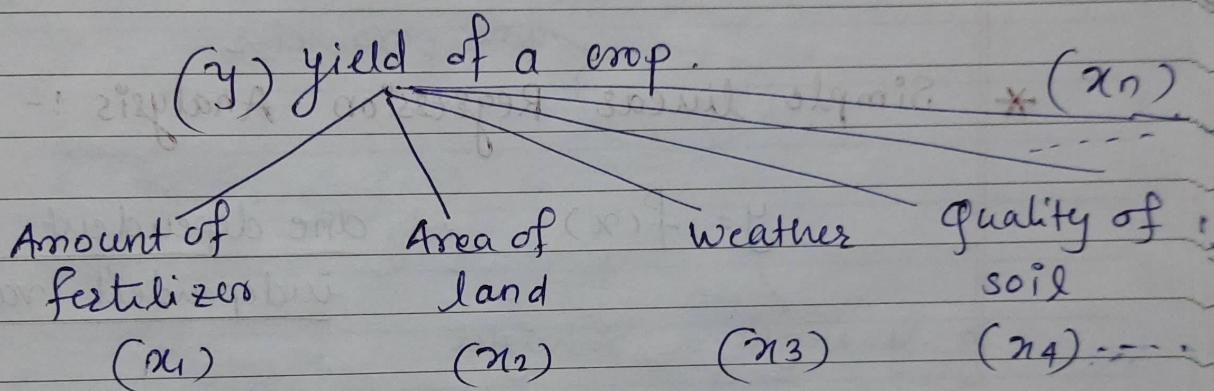
10th

Mondomery - D. C. Peck, E J Vining, G.G (2006)

Intro to linear Regression Analysis, Wiley.

* linear Regression is technique to model relationship
between ~~independent & dependent~~ & investigate

e.g.



∴ y depends on x_1, x_2, \dots, x_n . and all these are independent variables.

∴ There is a relation betw y and x_1, x_2, \dots, x_n .

$$\therefore y = f(x_1, x_2, \dots, x_n)$$

Aim \Rightarrow find out this f .

eg $y = x_1 + 2x_2 + 3x_3 + \dots + 10x_n$. linear
Relationship

eg $y = x_1x_2 + 2x_2^2 + 5x_3 + 6x_1x_4$ Non-linear
Relationship

Target is to find out Relationship based on data.

Model $\Rightarrow y = f(x_1, x_2, \dots, x_n)$

$x_1, x_2, \dots, x_n \Rightarrow$ Independent variables / Regressions

$y \Rightarrow$ Response variable / output var. / dependent variable.

* Simple linear Regression Analysis :-

$y = f(x)$ only one dependent & one independent variable.

e.g. $y = \text{Weight}$ $x = \text{height}$

(linear in terms of B_0 & B_1)

→ * linear wrt parameters, not variables.

$$\text{e.g. } y = B_0 + B_1 x^2$$

if $x^2 = Z$ then

$$y = B_0 + B_1 Z$$

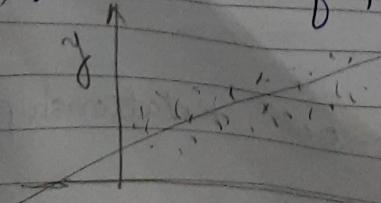
$$y = B_0 + B_1 x^2$$

$$y = B_0 + B_1 e^x$$

$$\text{e.g. } y = e^{B_0 + B_1 x}$$

Once you have data, you plot the data & if there is linear relationship then find values of B_0 & B_1 .

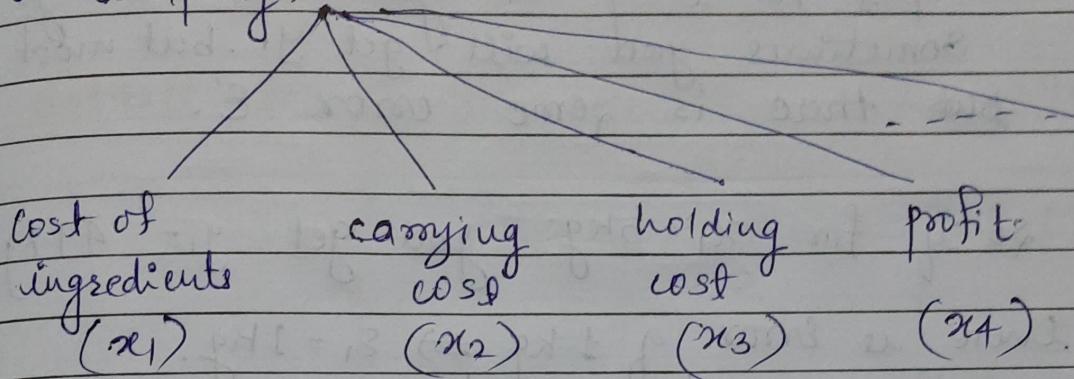
fitting
the model



Scatter plot
x

linear

eg. Price of an end product in a manufacturing company.



$$y = f(x_1, x_2, \dots, x_n)$$

e.g. x
↑
fertilizer y
↑
yield of crop

$$x_1 = 5 \text{ kg} \quad x_2 = 4 \text{ kg} \quad \dots \quad x_n = ? \text{ kg}$$

$$y_1 = 50 \text{ kg} \quad y_2 = 45 \text{ kg} \quad \dots \quad y_n = 60 \text{ kg}$$

∴ If we take avg. then we can say
on an avg. for 5 kg fertilizer, you are
getting 52 kg of crop.

$$\begin{array}{lll} x_1 = 5 & \rightarrow 50 \\ x_2 = 5 & \rightarrow 45 \\ x_3 = 5 & \rightarrow 55 \\ x_4 = 5 & \rightarrow 60 \end{array}$$

$$\therefore \frac{50+45+55+60}{4} = 52 \text{ kg}$$

on avg. for 5 kg

All the model might not fit

But for every x_i you don't always sometimes you will get y_i . But most time there is some error ' ϵ '.

so if for $x_1 = 5\text{ kg}$ you get $y_1 = 49\text{ kg}$ means there is error of 1 kg $\Rightarrow \epsilon_1 = 1\text{ kg}$.

$$E(y_i) \Rightarrow \beta_0 + \beta_1 x_1$$

so once you fit $x=5\text{ kg}$, on an avg how much yield you are getting that is expectn of y

Mathematic Model Vs Statistical Model

Vfor



Concrete



Actually takes care of what is happening

$$y_1 = \beta_0 + \beta_1 x_1$$

$$y_1 = \beta_0 + \beta_1 x_1 + \epsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_2$$

$$y_2 = \beta_0 + \beta_1 x_2 + \epsilon_2$$

$$y_n = \beta_0 + \beta_1 x_n$$

$$y_n = \beta_0 + \beta_1 x_n + \epsilon_n$$

i.e. If you fit $x_1 = 7\text{ kg}$ then on an avg you will get 60kg of crop

$$E(y_1) = \beta_0 + \beta_1 x_1, E(\epsilon_1) = 0$$

Unknown
QR Random
Error

M

data

The

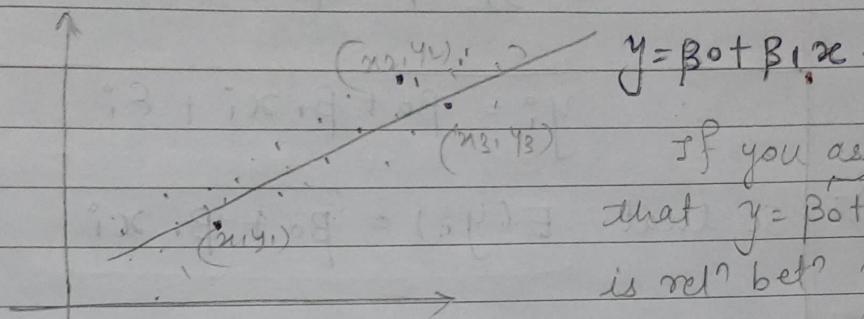
But

rat

Random
as
it is
RV

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

dataset $\Rightarrow (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$



The eqⁿ $y = \beta_0 + \beta_1 x$ should satisfy. Then all points should lie on a line. So

$$y_1 = \beta_0 + \beta_1 x_1$$

$$y_2 = \beta_0 + \beta_1 x_2$$

$$y_n = \beta_0 + \beta_1 x_n$$

Mathematical Model

But you don't actually get exactly response variable always rather, there is some error.

Random as of RV

$$\begin{aligned} y_1 &= \beta_0 + \beta_1 x_1 + \epsilon_1 \\ y_2 &= \beta_0 + \beta_1 x_2 + \epsilon_2 \\ &\vdots \\ y_n &= \beta_0 + \beta_1 x_n + \epsilon_n \end{aligned}$$

deterministic

Statistical Model

Unknown and is

Random Error (R.V.)

$$E(y_1) = \beta_0 + \beta_1 x_1, \quad E(\epsilon_1) = 0$$

$$E(y_n) = \beta_0 + \beta_1 x_n, \quad E(\epsilon_n) = 0$$

Once we know parameters
 β_0 and β_1 , will know the relationship betw x and y

Date _____
 Page _____

So our model in general is

dataset should
 satisfy this
 model

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

ϵ = error.

$$E(y) = \beta_0 + \beta_1 x, \text{ i.e., } E(\epsilon) = 0.$$

$$\therefore y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i=1, 2, \dots, n$$

JMP

$$\text{so that } E(y_i) = \beta_0 + \beta_1 x_i$$

$$\text{where } E(\epsilon_i) = 0, \quad \forall i = 1, 2, \dots, n$$

Assumptions \Rightarrow Observations:

① Random Errors ϵ_i 's are such that

$$E(\epsilon_i) = 0 \quad \text{and} \quad \text{var}(\epsilon_i) = \sigma^2 \quad \begin{matrix} \text{variability is} \\ \text{unknown f} \\ \text{fixed} \end{matrix}$$

② $\text{cov}(\epsilon_i, \epsilon_j) = 0, \quad \forall i \neq j$ Error of i th f j th
obs are not related
 i.e. ϵ_i & ϵ_j are unrelated.

③ $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$ Errors are iid RV
 f has $N(0, \sigma^2)$ distribution
 i.e. on avg error is zero

$$E(\epsilon_i) = 0 \quad \text{f. var}(\epsilon_i) = \sigma^2$$

f follows std normal distribution

variability
(σ^2)

std
normal distn
 $N(0, \sigma^2)$

Higher prob
of being on
line of reg.
prob.

target $y = \beta_0 + \beta_1 x_i + \epsilon$,

to find β_0 & β_1 and on avg you can have ϵ error.

* Least Square Estimation Method :-(Point estimation) :

find out β_0 and β_1 s.t. sum of the sq^r of error is min

$$S(\beta_0, \beta_1) = \sum_{i=1}^n \epsilon_i^2$$

is minimum (sum of squares of error is min)

i.e. all points should lie on the line of those theoretically most not lying on line should be close to line)

$$\epsilon_i = y_i - \beta_0 - \beta_1 x_i$$

fn of β_0, β_1
s.t. fn should take min value

Should take minimum value.

$$\therefore S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

This becomes problem of calculus now

$$\frac{\partial S}{\partial \beta_0} = 0 \quad \text{and} \quad \frac{\partial S}{\partial \beta_1} = 0, \text{ find critical points}$$

and then check Hessian matrix

d
sc

$$H = \begin{vmatrix} \frac{\partial^2 S}{\partial \beta_0^2} & \frac{\partial^2 S}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 S}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 S}{\partial \beta_1^2} \end{vmatrix}$$

If you do this minimization problem then,

2nd derivatives should come -ve &

H matrix should come +ve then you can say that ~~optimal~~ Local minimum and then check at Boundary pts f write conclusion

Global minimum based on two b.f.

Passage To find minimum of $S = 13 B = (1, 9, 10)$

will not go below during this time
and above will go beyond than
(will of 320)

$$1x.9 - 9 - 10 = 13$$

$$(1x.9 - 9 - 10)B = (1, 9, 10)$$

After Least Sq^t Estimate.

No. of Passengers (x)	Cost (\$1000) (y)
10	10
15	20
20	30
25	40
30	50
35	60
40	70
45	80
50	90
55	100
60	110
65	120
70	130
75	140
80	150
85	160
90	170
95	180
100	190

then:

$$\rightarrow \text{Re}l^n \Rightarrow \hat{y} = \beta_0 + \beta_1 x$$

$$\text{fitted model} \Rightarrow \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$$

(1) first calculate $\hat{\beta}_1$

$$(1) \hat{\beta}_1 = \frac{s_{xy}}{s_{xx}}$$

$$s_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x})$$

$$= \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)$$

$$① Sxx = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2$$

$$② Sxy = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)$$

$$③ \hat{\beta}_1 = \frac{Sxy}{Sxx}$$

$$④ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$⑤ \hat{\sigma}^2 = MSRes = \frac{SSRes}{n-2} = \frac{Syy - \hat{\beta}_1 Sxy}{n-2}$$

$$⑥ SSR = \hat{\beta}_1 Sxy$$

$$⑦ SST = SSRes + SSR$$

$$Syy$$

$$⑧ R^2 = \frac{SSR}{SST} = 1 - \frac{SSRes}{SST}$$

Least Square Estimates for $\hat{\beta}_0$ and $\hat{\beta}_1$

(Imp page)

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$$

Sum of
sqrsq nis

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\therefore S_{xy} \Rightarrow \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)$$

replace

y by x

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad , \quad \bar{y} = \frac{1}{n} \cdot \sum_{i=1}^n y_i$$

the fitted Regression model is given by.

$$\boxed{\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x}$$

fitted Regression
line using
least square Estimatⁿ
Method.

To check

(Point estimation)
How good you have fitted the model $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$
check properties for LSE for that.

CLASSMATE

Date _____

Page _____

18th Mon

(whether it is unbiased estimator of $MSE=0$)

Properties of LSE :- $\hat{\beta}_0$ & $\hat{\beta}_1$ are least square estimators.

$$1) \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

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

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

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\sum x_i = n\bar{x}$$

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

$$\boxed{\hat{\beta}_1 = \sum_{i=1}^n c_i y_i \text{, where } c_i = \frac{x_i - \bar{x}}{S_{xx}}}$$

So this estimator $\hat{\beta}_1$ is linear in terms of observation.

c_i satisfies following properties

① $\hat{\beta}_1$ is linear estimator

* ① $\sum_{i=1}^n c_i = 0$ ② sum of coefficient terms is zero.

* ② $\sum c_i x_i = 1$ ③ multiplying by independent variable of summing up will give value 1

"LSE is point estimation of pt estimation doesn't provide confidence or certainty."

* Point estimation \Rightarrow No guarantee.

for Least Sq^r check \rightarrow ① Unbiasedness

$$\text{② } \text{MSE} = \sigma^2$$

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i$$

We know $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

$$= \frac{1}{n} \sum_{i=1}^n y_i - \sum_{i=1}^n c_i y_i \bar{x}$$

$$\hat{\beta}_0 = \sum_{i=1}^n d_i y_i \quad \text{where } d_i = \frac{1}{n} - c_i \bar{x}$$

so $\hat{\beta}_0$ is also linear estimator of observation.

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i, \text{ where, } c_i = \frac{x_i - \bar{x}}{S_{xx}}$$

$$\hat{\beta}_0 = \sum_{i=1}^n d_i y_i, \text{ where } d_i = \frac{1}{n} - c_i \bar{x}$$

Both $\hat{\beta}_1$ and $\hat{\beta}_0$ are linear estimators.

$y = \beta_0 + \beta_1 x + \epsilon$ — starting model with error

↓ put the data & use least square estimation to get the fitted model

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ — fitted model.

$$E(\hat{\beta}_0) = \beta_0 \quad \text{and} \quad E(\hat{\beta}_1) = \beta_1$$

$\therefore \hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased estimators of β_0 & β_1

$$\text{var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad \begin{matrix} \text{How much } \hat{\beta}_1 \text{ can vary} \\ \text{or } S_{xx} - \text{sum of } x_i^2 \end{matrix}$$

$$\text{var}(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right] \quad \begin{matrix} \text{How much } \hat{\beta}_0 \text{ can vary} \\ \text{or } S_{xx} - \text{sum of } x_i^2 \end{matrix}$$

Gauss-Markov Thm :-

For the regression model $y = \beta_0 + \beta_1 x + \epsilon$,
with $E(\epsilon) = 0$ and $\text{var}(\epsilon) = \sigma^2$ and

uncorrelated error $[\text{cov}(\epsilon_i, \epsilon_j) = 0]$, the least squared estimators are unbiased as well as minimum variance.

Among all unbiased estimators, least sqr estimator method has min variance i.e., we get good line.
so we can say its a good estimator

#

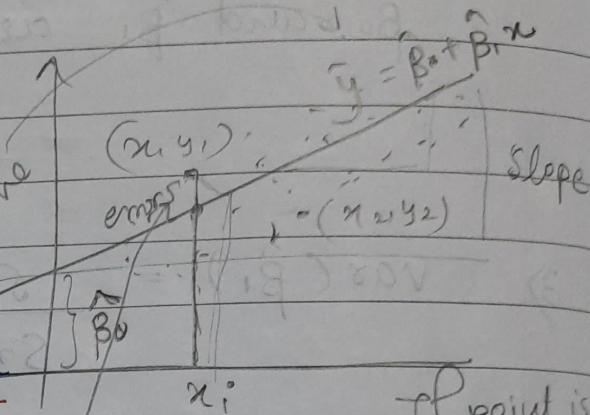
But how good you have fitted the line?

- 4) The sum of the residuals in any regression model that contains intercept term, is zero.

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

↓ ↓
intercept slope

Residual = obsv value - fitted value



If point is
on line then
Residual = 0

$$\text{Residual } (e_i) = y_i - \hat{y}_i$$

Then sum of the residuals is zero i.e.

$$\sum_{i=1}^n e_i = 0$$

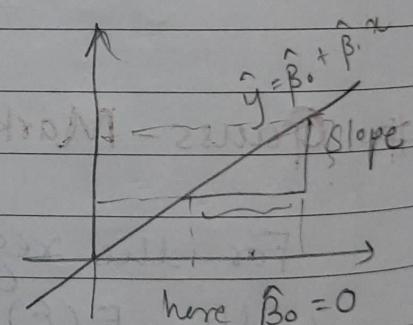
$$\Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i)$$

$$\Rightarrow \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)$$

$$\Rightarrow 0 \quad \frac{\partial S}{\partial \beta_0} = 0 \quad \text{check.}$$

$$S(\beta_0, \beta_1) = \sum_{i=1}^n \epsilon_i^2$$

$$\frac{\partial S}{\partial \beta_0} = 0 \Leftarrow y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i^0 \quad \text{Normal Eqn.}$$



here $\beta_0 = 0$

sum of the observed value is same as sum of fitted value.

$$\sum_{i=1}^n (y_i - \hat{y}_i) = 0$$

$$\boxed{\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i}$$

6) The least sq^r regression line passes through centroid (\bar{x}, \bar{y})

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ this eq^r of line is always passing through points (\bar{x}, \bar{y}) . which is centre of the data.

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$$

* All these properties are used to see how good our estimation is. i.e. least Sq^r estimator is one of the good estimator to estimate para of Regression model.

① Unbiasedness ✓

② Min Variance

③ sum of fitted value = sum of observed value

④ passes through origin of data

fitted model $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

* Estimation procedure for unknown variance σ^2

classmate
Date _____
Page _____
12th evening

Assumptions \Rightarrow

① $E(\varepsilon_i) = 0 \quad \text{and} \quad \text{var}(\varepsilon_i) = \sigma^2$ — fixed but unknown

② $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$

③ $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

Residual (e_i) = $(y_i - \hat{y}_i)$

Estimation of σ^2 :-

The sum of the square of the residual

$$SS_{\text{Res}} = \sum_{i=1}^n e_i^2$$

$$SS_{\text{Res}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Residuals

$$e_1 = y_1 - \hat{y}_1$$

$$\Rightarrow e_2 = y_2 - \hat{y}_2$$

$$SS_{\text{Res}} = Syy - \hat{\beta}_1 Sxy$$

// Random Variable

$$\hat{\beta}_1 = \frac{Sxy}{Sxx}$$

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

$$Sxx = \sum_{i=1}^n (x_i - \bar{x})^2$$

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

Result \Rightarrow

$$\frac{SS_{Res}}{\sigma^2} \sim \chi^2_{(n-2)}$$

$$\frac{(n-1)\sigma^2}{\sigma^2} \sim \chi^2_{n-1}$$

corollary \Rightarrow

$$X \sim \chi^2_{n-2}$$

$$\text{then } E(X) = n-2.$$

$$\frac{SS_{Res}}{\sigma^2} \sim \chi^2_{(n-2)}$$

$$X \sim \chi^2_{(n-2)}$$

$$\therefore E(X) = n-2.$$

$$\text{then } E\left(\frac{SS_{Res}}{\sigma^2}\right) = n-2$$

$$\Rightarrow \frac{1}{n-2} E(SS_{Res}) = \sigma^2$$

$$E(ax) = a E(x)$$

$E\left(\frac{SS_{Res}}{n-2}\right) = \sigma^2$
Mean Sq² Residual.

MSRes.

$$E(MSRes) = \sigma^2 \quad \text{where } MSRes = \frac{SS_{Res}}{n-2}$$

\therefore Estimate for σ^2
 Estimate for σ^2

$\hat{\sigma}^2 = MSRes$
estimate for σ^2

$$\frac{SS_{Res}}{n-2}$$

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

How good your model is
 # Evaluate Model \Rightarrow Hypothesis testing of the slope & intercept

$$\textcircled{1} \quad H_0: \beta_1 = 0 \\ \text{vs } H_1: \beta_1 \neq 0$$

$$\textcircled{2} \quad H_0: \beta_1 = 3 \\ \text{vs } H_1: \beta_1 \neq 3.$$

for intercept $\textcircled{3} \quad H_0: \beta_0 = \beta_0^* \\ \text{vs } H_1: \beta_0 \neq \beta_0^*$

If you accept your null or reject null then what will be conclusion?

Hypothesis testing problem \Rightarrow

13th morning.

* original model $\Rightarrow y = \beta_0 + \beta_1 x + \epsilon$

* fitted model $\Rightarrow \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

and assumptions are

- ① $E(\epsilon_i) = 0, \text{ var}(\epsilon_i) = \sigma^2$
- ② $\text{cov}(\epsilon_i, \epsilon_j) = 0$
- ③ $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

You want to test your belief, also checking how good your estimator is

The

eg $\beta_0 = \beta_1 = \beta_1^*$

Testing for Slope (β_1)

Testing problems for slope (β_1) parameter

$$H_0: \beta_1 = \beta_1^*$$

vs

$$H_1: \beta_1 \neq \beta_1^*$$

$$H_0: \beta_1 = \beta_1^*$$

vs

$$H_1: \beta_1 > \beta_1^*$$

$$H_0: \beta_1 = \beta_1^*$$

vs

$$H_1: \beta_1 < \beta_1^*$$

Testing Problems for Intercept also same as above.

* Result \Rightarrow

$$\text{i) } Z_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{s^2}{Sxx}}} \sim N(0,1) \quad \hat{\beta}_1 = \frac{Sxy}{Sxx}$$

$$\text{ii) } T_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{MSRes}{Sxx}}} \sim t(n-2)$$

* Testing Problems \Rightarrow

$$H_0: \beta_1 = \beta_1^* \quad \text{vs} \quad H_1: \beta_1 \neq \beta_1^*$$

① Testing slope (β_1) parameters:

case I $\Rightarrow \sigma^2$ is known

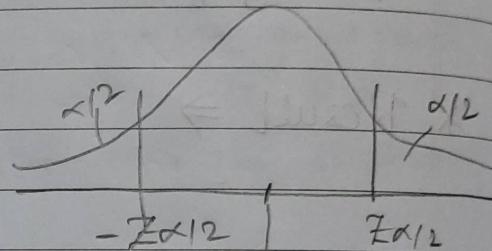
$$\textcircled{1} \quad H_0: \beta_1 = \beta_1^* \quad \text{vs} \quad H_1: \beta_1 \neq \beta_1^*$$

case I : σ^2 is known

$$Z_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{S_{xx}}}} \sim N(0,1)$$

level- α -test:

Reject H_0 if



$$|Z_1| = \left| \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{S_{xx}}}} \right| > Z_{\alpha/2}$$

↳ either $Z_1 > Z_{\alpha/2}$ or

$$\text{P-value} \Rightarrow 2P(Z_1 > |Z_{\text{obs}}|)$$

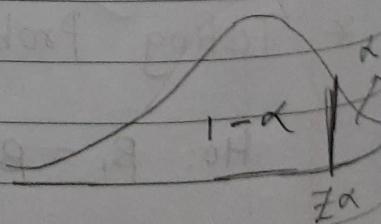
or

$$Z_1 < -Z_{\alpha/2}$$

$$\textcircled{2} \quad H_0: \beta_1 = \beta_1^* \quad \text{vs} \quad H_1: \beta_1 > \beta_1^*$$

Reject H_0 if

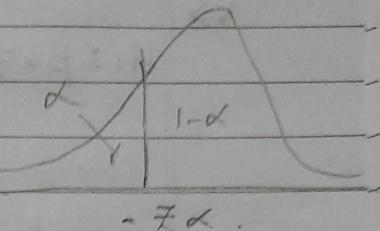
$$Z_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{S_{xx}}}} > Z_\alpha$$



$$\text{P-value} = P(Z_1 > Z_{\text{obs}})$$

(2) $H_0: \beta_1 = \beta_1^* \text{ vs } H_1: \beta_1 < \beta_1^*$

Reject H_0
if $Z_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{S_{xx}}}} < -Z_{\alpha}$



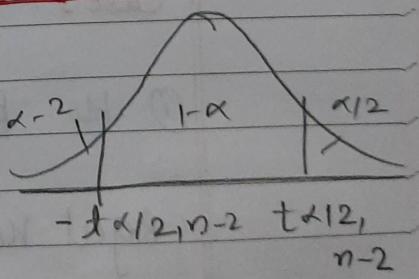
$$\text{P-value} = P(Z_1 < Z_{\text{obs}})$$

* If σ^2 is unknown, we have to use estimate for it
 $\hat{\sigma}^2 = \text{MSRes} = \frac{\text{SSRes}}{n-2} = \frac{S_{yy} - \hat{\beta}_1 S_{xx}}{n-2}$

Case II: σ^2 is unknown

(1) $H_0: \beta_1 = \beta_1^* \text{ vs } H_1: \beta_1 \neq \beta_1^*$

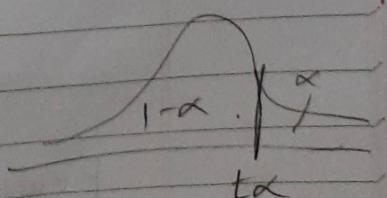
Reject H_0
if $|T_1| = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\text{MSRes}}{S_{xx}}}} > t_{(\alpha/2), n-2}$



$$\text{p-value} = 2 P(T_1 > |T_{\text{obs}}|)$$

(2) $H_0: \beta_1 = \beta_1^* \text{ vs } H_1: \beta_1 > \beta_1^*$

Reject H_0
if $T_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\text{MSRes}}{S_{xx}}}} > t_{\alpha, n-2}$



$$\text{---} < t_{\alpha, n-1}$$

② Testing intercept (β_0) parameter

$H_0: \beta_0 = \beta_0^*$	$H_0: \beta_0 = \beta_0^*$	$H_0: \beta_0 = \beta_0^*$
vs	vs	vs
$H_1: \beta_0 \neq \beta_0^*$	$H_1: \beta_0 > \beta_0^*$	$H_1: \beta_0 < \beta_0^*$

Results \Rightarrow

$$\textcircled{i} \quad Z_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim N(0,1)$$

σ^2 known \rightarrow

$$\textcircled{ii} \quad T_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MS_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim t(n-2)$$

σ^2 unknown \rightarrow

Case I: σ^2 known (same)

$$\textcircled{i} \quad H_0: \beta_0 = \beta_0^* \quad \text{vs} \quad H_1: \beta_0 \neq \beta_0^*$$

Reject H_0 if

$$|Z_0| = \left| \frac{\hat{\beta}_0 - \beta_0^*}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \right| > Z_{\alpha/2}$$

$$P\text{value} = 2 P(Z_0 > |Z_{0\text{obs}}|)$$

Same for
other
cases

classmate

Date _____

Page _____

② $H_0: \beta_0 = \beta_0^*$

vs $H_1: \beta_0 \neq \beta_0^*$

③ $H_0: \beta_0 = \beta_0^*$
vs $H_1: \beta_0 < \beta_0^*$

① $H_0:$

vs $H_1:$

② $H_0: \beta_0 = \beta_0^*$

vs $H_1: \beta_0 \neq \beta_0^*$

③ $H_0: \beta_0 = \beta_0^*$

vs

$H_1: \beta_0 < \beta_0^*$

classmate
Date _____
Page _____

Reject Ho if

① $H_0: \beta_0 = \beta_0^*$ $Z_0 = \frac{\hat{\beta}_0 - \beta_0^*}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} > Z_\alpha$

vs $H_1: \beta_0 > \beta_0^*$

③ $H_0: \beta_0 = \beta_0^*$ $Z_0 = \frac{\hat{\beta}_0 - \beta_0^*}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} < -Z_\alpha$

vs $H_1: \beta_0 < \beta_0^*$

Case II: σ^2 is unknown

① $H_0: \beta_0 = \beta_0^*$ Reject Ho if

vs $H_1: \beta_0 \neq \beta_0^*$

p-value = $2P(T_0 > |T_{obs}|)$

$$|T_0| = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MS_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \xrightarrow{\text{if } \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MS_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} > t_{\alpha/2, n-2}}$$

② $H_0: \beta_0 = \beta_0^*$

vs $H_1: \beta_0 > \beta_0^*$

$$T_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MS_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} > t_{\alpha, n-2}$$

③ $H_0: \beta_0 = \beta_0^*$

vs

$H_1: \beta_0 < \beta_0^*$

$$T_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MS_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} < -t_{\alpha, n-2}$$

$$\beta_1^* = 0$$

* $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

Based on data if you accept Null (H_0) means slope of the para is zero.

fitted model $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

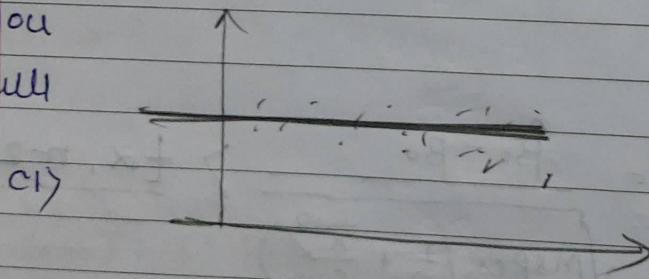
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

if we accept null then $\hat{\beta}_1 = 0$ slope zero.

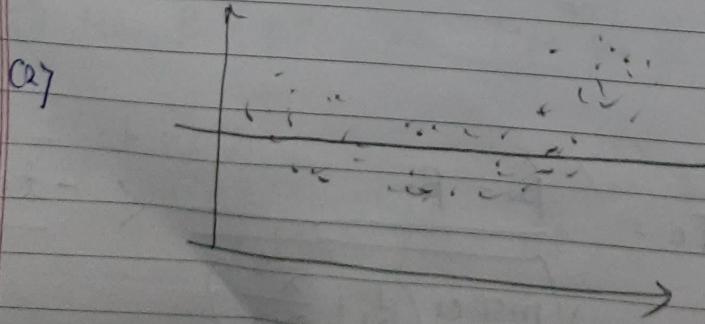
$$\hat{\beta}_0 = \bar{y}$$

$\therefore \hat{y} = \bar{y}$ for all the points x
fitted line will be
some constant

When you
Accept null



c) Regression line
parallel to x-axis

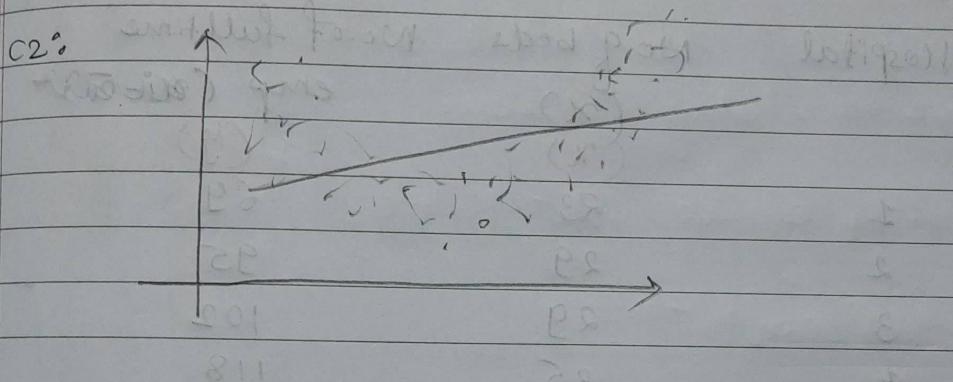
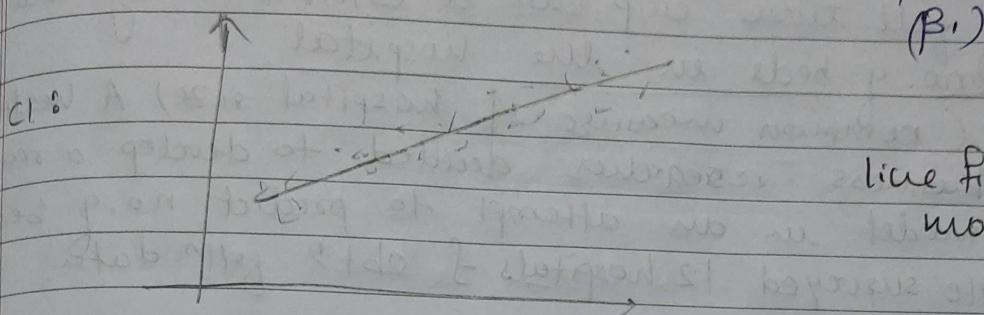


c) data has
non-linear
behaviour if
you are trying
to fit straight
line

13th evening

when you reject Null means slope para $\neq 0$

(B.)

non-linear
data

when you accept alternative either you fitted well the model OR the model doesn't fit well but if try to fit polynomial model it might fit well

$$y = \hat{\beta}_0 + \hat{\beta}_1 x^2 + \epsilon$$

13th evening

Ex. 1) A specialist in hospital admin stated that no. of full time emp can be estimated by counting no. of beds in the hospital

(common measure of hospital size) A healthcare business researcher decided to develop a regression model in an attempt to predict no. of beds. He surveyed 12 hospitals & obt'd foll' data.

Hospital	No. of beds (x)	No. of full time emp (dependent (y))
1	23	69
2	29	95
3	29	102
4	35	118
5	42	126
6	46	138
7	50	178
8	54	156
9	64	184
10	66	176
11	78	225
12	76	125

↑ dependent
variable

→ Once there is linear relationship then we fits the straight line

Least sqs estimate ⇒

fitted Regression

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad \text{--- } ①$$

classmate

Date _____

Page _____

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad \text{--- } ②$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \text{--- } ③$$

putting ③ in ①

$$\boxed{\hat{y} = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x \Rightarrow \hat{y} = \bar{y} - \hat{\beta}_1 (x - \bar{x})}$$

~~$$\hat{y} = \bar{y} - \hat{\beta}_1 (\bar{x} + x)$$~~

$$① \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

$$\sum x^o = 592, \quad \sum y^o = 1692.$$

$$\bar{x} = \frac{1}{n} \sum x^o$$

$$\bar{y} = \frac{1}{n} \sum y^o$$

$$\bar{x} = \frac{1}{12} \times 592$$

$$\bar{y} = \frac{1}{12} \times 1692$$

$$\boxed{\bar{x} = 49.33}$$

$$\boxed{\bar{y} = 141}$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = 3838.7$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 8566.$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{8566}{3838.7} = 2.232$$

$$\hat{y} = \bar{y} - \hat{\beta}_1 (x - \bar{x})$$

$$\hat{y} = 141 - 2.232(x - 49.33)$$

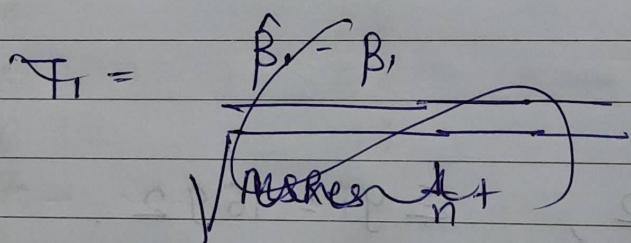
fitted eqn
of line.

If we have testing problem -

Case I $\Rightarrow H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

* Generally σ^2 are not known.

σ^2 unknown



Reject $|T_1| = \left| \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\frac{MSRes}{S_{xx}}}} \right| > t_{\alpha/2, n-2}$

$$T_1 = \frac{2.232 - 0}{\sqrt{\frac{244.886}{3838.7}}} = 8.84$$

$$MS_{Res} = \frac{SS_{Res}}{n-2} = \frac{s_{y-y-\hat{\beta}_1 x}}{n-2}$$

$$= \frac{2448.86 - 2.232 \times 8566}{11}$$

$$MS_{Res} = \frac{SS_{Res}}{n-2} = \frac{s_{y-y-\hat{\beta}_1 x}}{n-2}$$

$$\therefore MS_{Res} = 244.86.$$

$$\textcircled{2} \quad t_{\alpha/2, n-2} \Rightarrow$$

$$\textcircled{3} \quad p\text{-value} = 2P(T_1 > |T_{obs}|)$$

$$= 2P(T_1 > 8.84)$$

p-value = ~~2*~~ close to zero

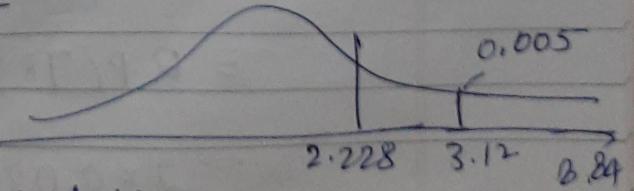
$$dof = n-2$$

$$= 10$$

If p-value $\leq \alpha$ then reject

for any α p-value is
close to zero so we reject H_0 .

the slope means β_1 is non-zero quantity.



α will be
close to zero

Case II: $H_0: \beta_0 = 0$ vs $H_1: \beta_0 \neq 0$ σ^2 unknown
(natural scenario)

$$T_0 = \frac{\hat{\beta}_0 - \beta_0^*}{\sqrt{MSRes \left(\frac{1}{n} + \frac{\bar{x}^2}{Sx^2} \right)}}$$

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ &= 141 - 2.232 \times 49.33\end{aligned}$$

$$\hat{\beta}_0 = 30.895$$

$$\beta_0^* = 0$$

$$MSRes = 244.886$$

$$Sx^2 = 3838.7$$

$$|T_0| = 2.333 > t_{\alpha/2, n-2} ?$$

$$P\text{-value} = 2P(T_0 > t_{\text{obs}})$$

$$= 2P(T_0 > 2.33) \quad \text{at } 10 \text{ d.o.f}$$

$$= 2 \times 0.035$$

$$= 0.07$$

$$\alpha = 0.035$$

If $\alpha = 0.05$

$$t_{\alpha/2, n-2} = t_{0.025, 10} = 2.228.$$

$$\therefore T_0 = 2.333 > t_{\alpha/2, n-2} = 2.228$$

so we reject Null. means intercept term $\neq 0$.

If $\alpha = 0.01$

$$t_{\alpha/2, n-2} = t_{0.005, 10} = 3.169.$$

$$T_0 = 2.333 \not> t_{\alpha/2, n-2} = 3.169$$

so we accept Null (H_0) i.e. -

$$MS_{Res} = \frac{SS_{Res}}{n-2}$$

classmate

Date _____

Page _____

10th M

* Interval Estimation in Simple linear Regression Analysis \Rightarrow

$$y = \beta_0 + \beta_1 x + \varepsilon. \quad \text{our model.}$$

Results $\Rightarrow \sigma^2$ unknown

$$\textcircled{i} \quad T_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{MS_{Res}}{S_{xx}}}} \sim t(n-2)$$

$$\textcircled{ii} \quad T_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MS_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim t(n-2)$$

$$MS_{Res} = \frac{SS_{Res}}{n-2}$$

Let

Estimated
Standard
Error

$$Se(\hat{\beta}_1) = \sqrt{\frac{MS_{Res}}{S_{xx}}} = \frac{s_{y\bar{x}} - \hat{\beta}_1 s_{x\bar{x}}}{n-2}$$

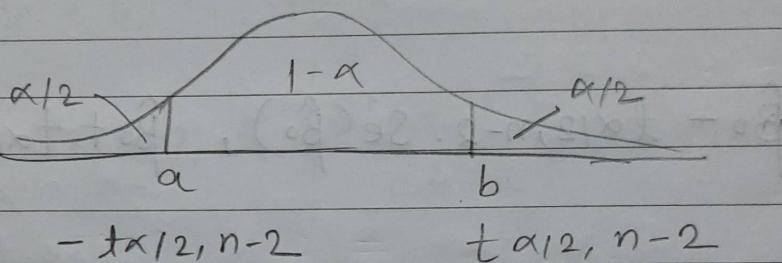
$$Se(\hat{\beta}_0) = \sqrt{MS_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}$$

$$\therefore T_1 = \frac{\hat{\beta}_1 - \beta_1}{Se(\hat{\beta}_1)} \quad \text{and} \quad T_0 = \frac{\hat{\beta}_0 - \beta_0}{Se(\hat{\beta}_0)}$$

These T_1 and T_0 will be used as pivot quantity to estimate parameters $\hat{\beta}_0$ & $\hat{\beta}_1$.

* Interval estimation for $\hat{\beta}_1 \xrightarrow{\beta_1}$ σ^2 unknown.

$100(1-\alpha)\%$ CI for parameter β_1 based on pivot quantity T_1 is given by,



$$\therefore P(a < T_1 < b) = 1-\alpha$$

$$P\left(-t_{\alpha/2, n-2} < \frac{\hat{\beta}_1 - \beta_1}{Se(\hat{\beta}_1)} < t_{\alpha/2, n-2}\right) = 1-\alpha$$

$$P\left(-t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_1) < \hat{\beta}_1 - \beta_1 < t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_1)\right) = 1-\alpha$$

$$P\left(\hat{\beta}_1 - t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_1) < \beta_1 < \hat{\beta}_1 + t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_1)\right)$$

$$\text{Interval} \Rightarrow (\hat{\beta}_1 - t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_1), \hat{\beta}_1 + t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_1))$$

* Interval Estimation for $\beta_0 \Rightarrow$

$$T_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MSE_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{Sxx} \right)}} \sim t(n-2)$$

100(1- α)% CI for β_0 based on T_0 is given by

$$\left(\hat{\beta}_0 - t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_0), \hat{\beta}_0 + t_{\alpha/2, n-2} \cdot Se(\hat{\beta}_0) \right)$$

where, $Se(\hat{\beta}_0) = \sqrt{MSE_{Res} \left(\frac{1}{n} + \frac{\bar{x}^2}{Sxx} \right)}$

* Estimation for model parameter σ^2

$$Var(\varepsilon_i) = \sigma^2 \quad \text{variance quantity of error}$$

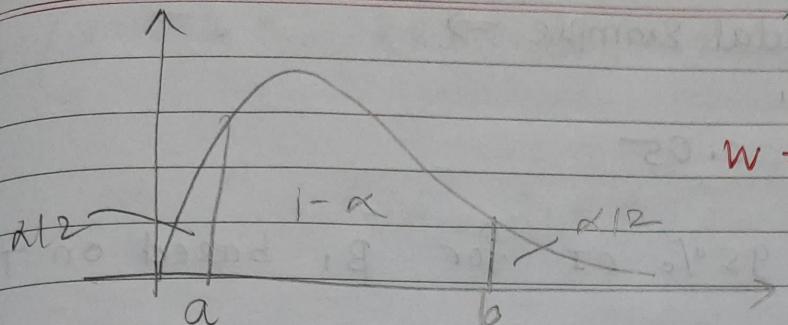
Result \Rightarrow

$$W = \frac{SS_{Res}}{\sigma^2} \sim \chi^2(n-2)$$

Interval Estimation of $\sigma^2 \Rightarrow$

classmate

Date _____
Page _____



$$W = \frac{SS_{\text{Res}}}{\sigma^2} \sim \chi^2_{n-2}$$

$$\chi^2_{1-\alpha/2, n-2} \quad \chi^2_{\alpha/2, n-2}$$

$$P(a < W < b) = 1 - \alpha.$$

$$P\left(\chi^2_{1-\alpha/2, n-2} < \frac{SS_{\text{Res}}}{\sigma^2} < \chi^2_{\alpha/2, n-2}\right) = 1 - \alpha.$$

$$P\left(\frac{SS_{\text{Res}}}{\chi^2_{\alpha/2, n-2}} < \sigma^2 < \frac{SS_{\text{Res}}}{\chi^2_{1-\alpha/2, n-2}}\right) = 1 - \alpha.$$

100(1 - α)% CI for σ^2 based on pivot quantity
W is

$$\left(\frac{SS_{\text{Res}}}{\chi^2_{\alpha/2, n-2}}, \frac{SS_{\text{Res}}}{\chi^2_{1-\alpha/2, n-2}} \right)$$

where, $SS_{\text{Res}} = Syy - \hat{\beta}_1 Sxy$.

Interval estimations for $\beta_1, \beta_0, \sigma^2 \Rightarrow$

Same hospital example \Rightarrow

① for β_1

Let $\alpha = 0.05$

$100(1-\alpha) = 95\%$ CI for β_1 based on T, is given by

$$(\hat{\beta}_1 - t_{\alpha/2, n-2} S_e(\hat{\beta}_1), \hat{\beta}_1 + t_{\alpha/2, n-2} S_e(\hat{\beta}_1))$$

$$\hat{\beta}_1 = 2.232$$

$$S_e(\hat{\beta}_1) = \sqrt{\frac{MSE}{S_{xx}}} = \sqrt{\frac{244.886}{3838.7}} =$$

$$= \sqrt{\frac{SS_{Res}}{(n-2)S_{xx}}}$$

$$= \sqrt{\frac{2448.86}{10 \times 3838.7}} \quad MSE = 244.8$$

$$S_e = 0.25$$

$$S_e(\hat{\beta}_1) = 0.25$$

$$t_{\alpha/2, n-2} = t_{0.025, 10} = 2.228$$

$$\therefore (2.232 - 2.228 \times 0.25, 2.232 + 2.228 \times 0.25)$$

$$\Rightarrow (1.675, 2.789)$$

β_0 in pt estimation $\hat{\beta}_1 = 2.232$.

in interval estimation we can say β_1 lies in interval $(1.675, 2.789)$ with 95% confidence.

② for β_0

$$\alpha = 0.05, \beta_0 = 30.895$$

$$(\hat{\beta}_0 - t_{\alpha/2, n-2} s_e(\hat{\beta}_0), \hat{\beta}_0 + t_{\alpha/2, n-2} s_e(\hat{\beta}_0))$$

$$s_e(\hat{\beta}_0) = \sqrt{MSE \left(\frac{1}{n} + \frac{\bar{x}^2}{Sxx} \right)}$$

$$s_e(\hat{\beta}_0) = 13.25$$

$$(1.368, 33.944) \text{ interval for } \beta_0$$

③ for f^2 - 95% CI for f^2 .

$$\left(\frac{\text{SSRes}}{\chi^2_{\alpha/2, n-2}}, \frac{\text{SSRes}}{\chi^2_{1-\alpha/2, n-2}} \right)$$

$$\text{At. } \alpha = 0.05$$

$$\text{SSRes} = 2448.6$$

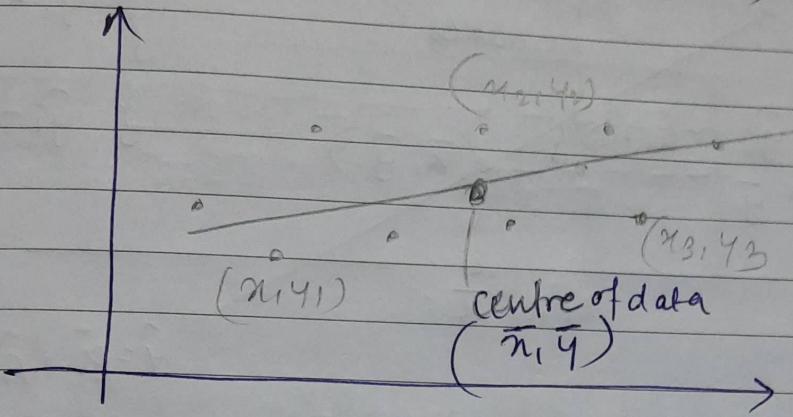
$$\chi^2_{\alpha/2, n-2} = \chi^2_{0.025, 10} = 20.483$$

$$\chi^2_{1-\alpha/2, n-2} = \chi^2_{0.975, 10} = 3.247$$

Interval for f^2 $(119.58, 753.49)$

* Analysis of Variance (ANOVA) \Rightarrow

$$(x_1, y_1) \quad (x_2, y_2) \quad \dots \quad (x_n, y_n)$$



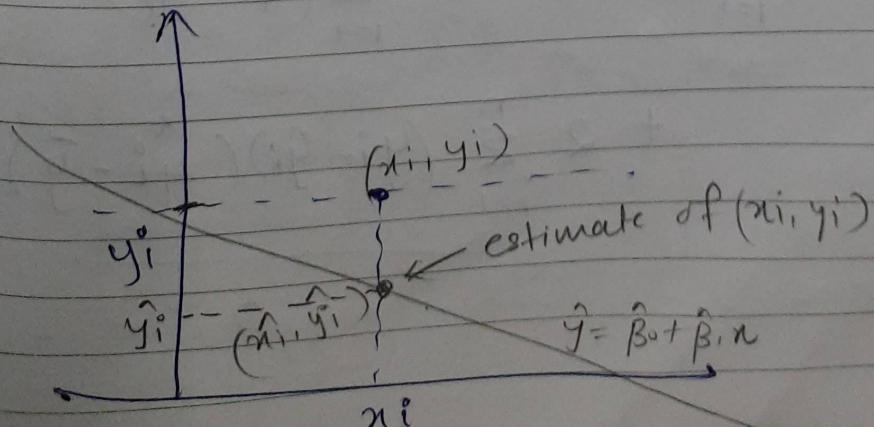
Avg Variability in data.

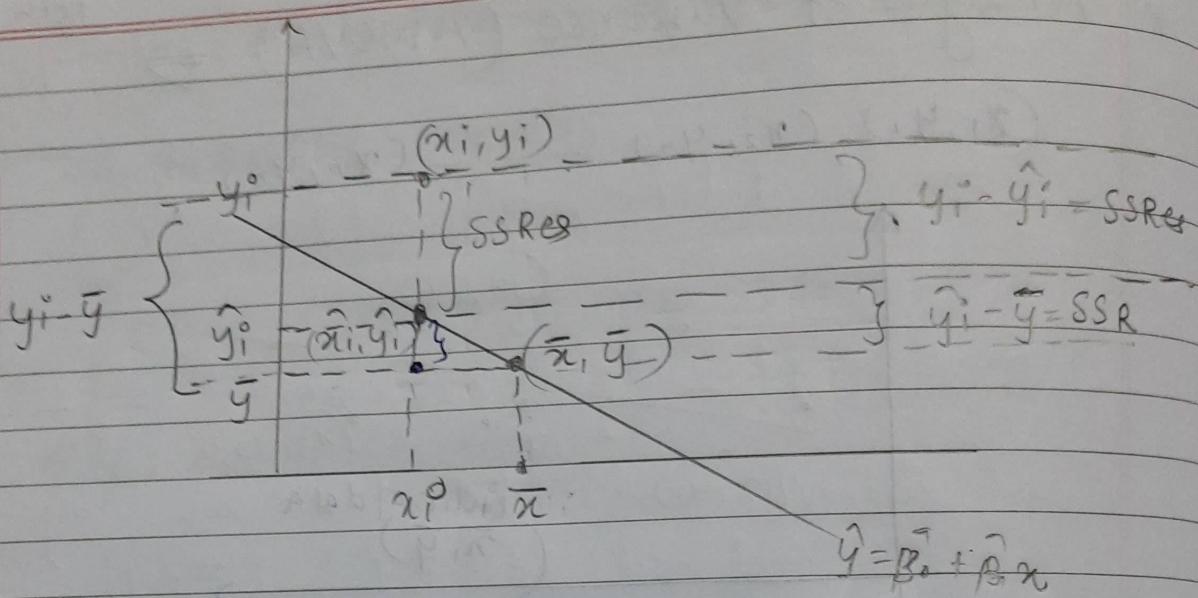
$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

how far
are my data
pts from
the centre
of the data

Any Regression line always passes
through centre of the data
centroid

$$\text{Total variability} = \sum_{i=1}^n (y_i - \bar{y})^2$$





$(y_i - \bar{y})$ is total variability in i th obsv.

deviation, $y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$ variation of data when x is not considered.
 $=$ ~~all~~ y_i

Total variability of whole data

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n [(y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})]^2$$

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

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

Sum of the squares of Total Variation

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

$$(SS_T \equiv Syy)$$

↓ SSRes

Sum of the
squares of
residual.

↓ SSR

Sum of the
square due
to regression

$$\therefore SS_T = SS_{\text{Res}} + SSR$$

$$\textcircled{1} \quad SS_{\text{Res}} = Syy - \hat{\beta}_1 Sxy$$

$$SS_{\text{Res.}} = SS_T - \hat{\beta}_1 Sxy$$

$$= Ssy -$$

$$SSR = SS_T - SS_{\text{Res}}$$

$$= SS_T - [SS_T - \hat{\beta}_1 Sxy]$$

$$\textcircled{2} \quad SSR = \hat{\beta}_1 Sxy$$

$$SS_{\text{Res}} = SS_T - \hat{\beta}_1 Sxy$$

← same

* Coefficient of determination \Rightarrow

The quantity

coeff of
determinat

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SS_{Res}}{SST}$$

how good
your model
fits.

Observations :-

① $R^2 = 1$ iff $SS_{Res} = 0$ (when sum of
sqr of residuals
is zero)

when all pts lie on fitted straight line.

(when obs is exactly on the line) $\sum_{i=1}^n (y_i - \hat{y}_i) = 0$

then Residual term = 0

so you fitted the model very well.

② $R^2 = \frac{SSR}{SST} \Rightarrow$ what proportion of ~~measur~~
variability is explained by
the data.

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

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

variability of y by
taking care of regression
variables.

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

variation of y 's

③ $R^2 = 0$, $SSR = 0 \rightarrow$ (when $\hat{y} = \bar{y}$)

or $SSRes = SST$.

④ $0 \leq R^2 \leq 1$ how good your model is fitted

y R^2 closer to 1 - Good fit

$R^2 \rightarrow 0$ - poor fit

$$SSR \leq SST$$

$$\frac{SSR}{SST} \leq 1$$

Example continued. \Rightarrow

$$Syy = 856.6$$

$$SSRes = 2448.86.$$

$$\hat{\beta}_1 = 2.232$$

$$\hat{y} = 141 + 2.232(x - 49.33)$$

$$SST = Syy = 211564$$

$$R^2 = 1 - \frac{SSRes}{SST} \Rightarrow 1 - \frac{2448.86}{211564} \Rightarrow 0.886$$

→ fitted line explains 88.6% variability of data if it cannot explain remaining 11.4% variability in data

$$⑤ R^2 \leq 1 \quad \hat{y} = \beta_0 + \beta_1 x + x^2$$

$$⑥ E(R^2) = 1 - \frac{\sigma^2}{\hat{\beta}_1 S_{xx} + \sigma^2}$$

$$S_{xx} \uparrow = E(R^2) \uparrow$$

$$S_{xx} \downarrow = E(R^2) \downarrow$$

* Analysis of Variance (ANOVA) \Rightarrow F-test

① F-distribution \Rightarrow

Let X and Y be two independent RV s.t.

$$X \sim \chi^2(n) \text{ and } Y \sim \chi^2(m)$$

$$F = \frac{X/n}{Y/m} \sim F_{n,m} \quad (\text{ratio of two chi-sqr dist's})$$

② Non-central Chi-sqr distribution \Rightarrow

$$x_i \stackrel{\text{iid}}{\sim} N(0,1)$$

$$\begin{array}{l} \text{Central } W = \sum_{i=1}^n x_i^2 \sim \chi^2(n) \\ \text{Chi-sqr} \end{array} \leftarrow \mu_i = 0, \forall i$$

Now if

$$x_i \stackrel{\text{independent}}{\sim} N(\mu_i, 1)$$

$$\begin{array}{l} \rightarrow W = \sum_{i=1}^n x_i^2 \sim \chi^2(n) \text{ with non-central para. } \lambda \\ \text{Non-central chi-sqr} \end{array}$$

$$\lambda = \sum_{i=1}^n \mu_i^2$$

Result \Rightarrow

$$\text{i) } \frac{SS_{\text{Res}}}{\sigma^2} \sim \chi^2(n-2)$$

\leftarrow usual chisq $\lambda = 0$

$$\text{ii) } \frac{SS_R}{\sigma^2} \sim \chi^2(1)$$

\leftarrow Non-central chisq
with non-central
para $\lambda = \hat{\beta}_1^2 S_{xx}$

$$\text{iii) } F_0 = \frac{MS_R}{MS_{\text{Res}}} \sim F_1, n-2$$

$$MS_R = \frac{SS_R}{dof} = \frac{SS_R}{1} \text{ mean sqr due to regression}$$

$$MS_{\text{Res}} = \frac{SS_{\text{Res}}}{n-2} \text{ mean sqr residual.}$$

$$MS_R = SS_R = \hat{\beta}_1 S_{xy}$$

original
model

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

fitted
model

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

* For simple linear regression T-test & F-test will give same results.

F-test \Rightarrow

$H_0: \beta_1 = \beta_1^*$ vs $H_1:$

$$\textcircled{1} \quad H_0: \beta_1 = 0 \quad \text{vs} \quad H_1: \beta_1 \neq 0$$

level- α -test \Rightarrow

Reject H_0 if

$$F_0 = \frac{MSR}{MS_{Res}} \rightarrow F_{\alpha, 1, n-2}$$

table F_{α, V_1, V_2}

* ANOVA Table \Rightarrow

Source of Variation	Sum of Squares	DOF	Mean Square	F ₀
---------------------	----------------	-----	-------------	----------------

① Regression $SSR = \hat{\beta}_1 S_{xy}$ (1) $MSR = \frac{SSR}{1}$

$$F_0 = \frac{MSR}{MSRes}$$

② Residual $SSRes = SST - \hat{\beta}_1 S_{xy}$ (n-2) $MSRes = \frac{SSRes}{n-2}$

Total $SST = SSR + SSRes$ n-1

Ex.1) Hospital Example \Rightarrow

ANOVA table \Rightarrow

$$SSR = \hat{\beta}_1 S_{xy} = 2.232 \times 3838.7 = 19119.3$$

$$SSRes = S_{yy} - \hat{\beta}_1 S_{xy} = 2444.86$$

$$MSR = SSR = 19119.3$$

Source of variation	Sum of Squares	D.F	Mean Sq	F.O
Regression	SSR = 19119.3	1	MSR = 19119.3	$F_O = \frac{19119.3}{244.46}$
Residual	SSRes = 24448.6	10	MSRes = 244.446	= 78.973

If $\alpha = 0.05$

$$F_{\alpha, 1, 10} \Rightarrow F_{0.05, 1, 10} \Rightarrow 4.96.$$

$$\therefore F_O = 78.973 > F_{\alpha, 1, 10} = 4.96.$$

∴ Reject H_0 i.e. $\beta_1 \neq 0$

Result \Rightarrow Relation betw T-test and F-test \Rightarrow

$$X \sim t(n) \text{ then } X^2 \sim F_{1, n}$$

dof

$$T_1 = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\frac{MSRes}{S_{xx}}}} \sim t(n-2)$$

$$T_1 = 8.84$$

$$\text{Jmp} \rightarrow T_1^2 = (8.84)^2 = 78.973 = F_0$$

$$T_{\text{obs}} > t_{\alpha/2, n-2}$$

If $\alpha = 0.05$

$$t_{0.025, 10} = 2.228.$$

$$T_1 = 8.84 > t_{\alpha/2, n-2} = 2.228$$

∴ We reject H_0

$$H_0: \beta_1 = 0 \text{ vs } H_a: \beta_1 \neq 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{F-test not applicable.}$$

$$\text{OR } H_0: \beta_1 = 0 \text{ vs } H_a: \beta_1 > 0 \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

F-test is applicable only for two sided hypothesis.

* Application of Simple linear Regression Analysis \Rightarrow

corresponding to a regression variable, how do you predict the mean response value.

$$\text{our model} \Rightarrow y = \beta_0 + \beta_1 x + \epsilon$$

$$\text{fitted model} \Rightarrow \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

If $x = x_0$ (if $x = 40$ beds in hospital) then
 what is $E(y/x=x_0)$ (no of full time emp
 on an avg
 when $x=40$ beds)

$$E(\hat{y}/x=x_0) \equiv \hat{M}y/x_0$$

① Point Estimation \Rightarrow

$$y = \beta_0 + \beta_1 x_0 + \epsilon$$

$$\hat{M}y/x_0 = E(\hat{y}/x=x_0) = \beta_0 + \beta_1 x_0 \quad (\because E(\epsilon)=0)$$

$$\hat{M}y/x_0 = E(y/x=x_0) = \beta_0 + \beta_1 x_0 \quad (\because E(\epsilon)=0)$$

$$\boxed{\hat{M}y/x_0 = E(\hat{y}/x=x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0}$$

point of
 estimation
 of mean
 response of y
 given $x=x_0$

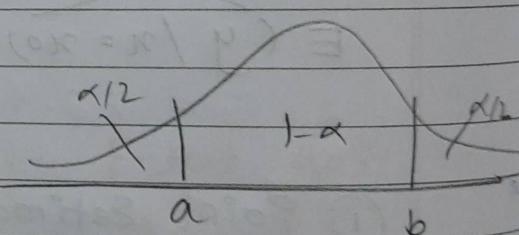
② Interval Estimation

If $n=40$ beds then what is 99% chance that or on an avg there will be how many no of employees.

Result \Rightarrow

$$\text{pivot, } T = \frac{\hat{U}_{Y/x_0} - U_{Y/x_0}}{\sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{Sxx} \right)}} \sim t(n-2)$$

$$P(a < T < b) = 1 - \alpha$$



$$P\left(-t_{\alpha/2, n-2} < \frac{\hat{U}_{Y/x_0} - U_{Y/x_0}}{\sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{Sxx} \right)}} < t_{\alpha/2, n-2}\right) = 1 - \alpha$$

$$P\left(\hat{U}_{Y/x_0} - t_{\alpha/2, n-2} \sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{Sxx} \right)} < U_{Y/x_0}\right)$$

$$\left\langle \hat{U}_{Y/x_0} + t_{\alpha/2, n-2} \sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{Sxx} \right)} \right\rangle = 1 - \alpha$$

$100(1-\alpha)\%$ PCT for \hat{y}_{y/x_0} based on t is

$$\left(\hat{y}_{y/x_0} - t_{\alpha/2, n-2} \sqrt{MS_{Res}} \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_x^2} \right) , \right.$$

$$\left. \hat{y}_{y/x_0} + t_{\alpha/2, n-2} \sqrt{MS_{Res}} \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_x^2} \right) \right)$$

Prediction of Mean Response corresponding to
a given value $x = x_0$

$$\hat{y}_{y/x_0} = E(\hat{y}|x=x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

Example continued.

Q. If there are 40 beds then how many full time employees?

→ ① point estimate

$$\hat{y}_{y/x_0} = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

$$= 30.895 + 2.322 \times 40$$

$$= 120.17 \cdot 120 \text{ full time employees}$$

② Interval estimation \Rightarrow

95% CI then $\alpha = 0.05$

$$t_{\alpha/2, n-2} = t_{0.025, 10} = 2.228$$

$$n_0 = 40, \bar{x} = 49.33, \bar{y} = 141, S_{xy} = 3838.7,$$

$$\hat{\beta}_1 = 2.232, \hat{\beta}_0 = 30.895$$

$$MS_{Res} = 244.886$$

$$\hat{y}_{x_0} = 120.17$$

find values

$$(108.82, 131.52)$$

$$\approx (109, 132) \text{ no of full time employees}$$

* Observations

classmate

Date _____

Page _____

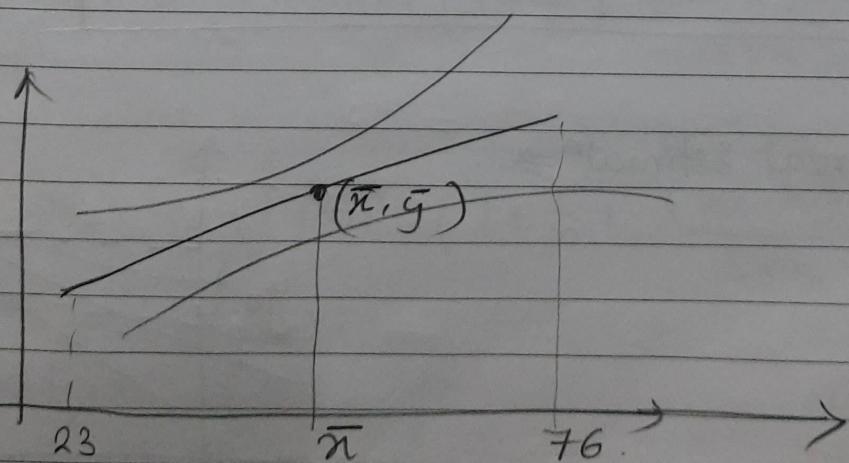
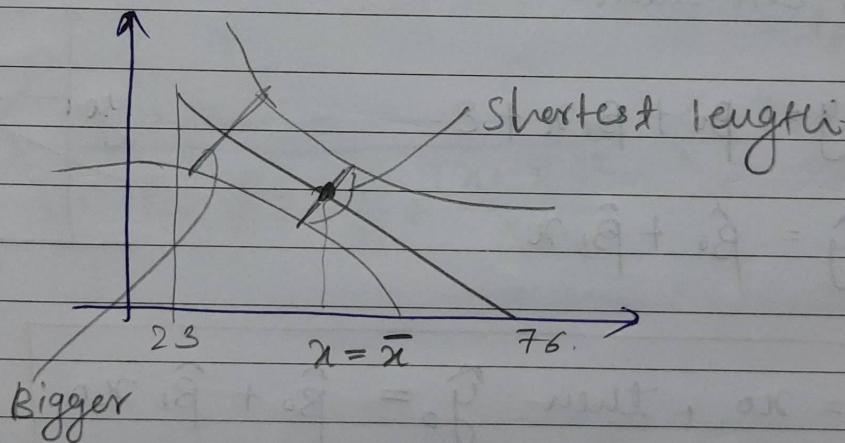
The length of CI is \overrightarrow{ab}

$$a(b-a) = 2 \times 12, n-2 \sqrt{MSE} \left(\frac{1}{n} + \frac{(n_0 - \bar{n})^2}{3n^2} \right)$$

$$\therefore (b-a) =$$

~~length~~ Then length of CI will be minimum when

$$n_0 = \bar{n}$$



for some point outside the range, there is chance of more error.

Always suggested to do ~~base~~ prediction for
the given value.

prediction of New Observation \Rightarrow

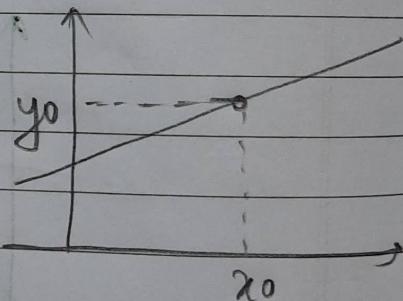
predict future obs.

If $x = x_0$ then predict $y = y_0$? (not asking mean $E(y/x=x_0)$)

① point estimation \Rightarrow

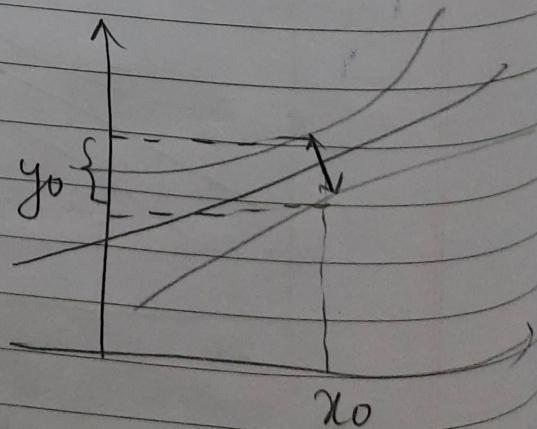
$$y = \beta_0 + \beta_1 x + \varepsilon.$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$



$$x = x_0, \text{ then } \hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

② Interval Estimation \Rightarrow

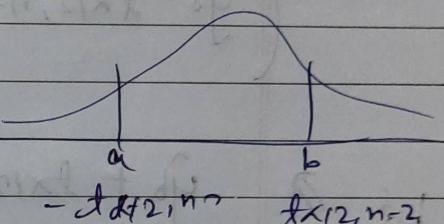


Result \Rightarrow

$$T^* = \frac{\hat{y}_0 - y_0}{\sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}} \sim t(n-2)$$

$100(1-\alpha)\%$ CI based on T^* is given as.

$$P(a < T^* < b) = 1 - \alpha$$



$$P(-t_{\alpha/2, n-2} < \frac{\hat{y}_0 - y_0}{\sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}} < t_{\alpha/2, n-2})$$

$$P(\hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} < y_0 < \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)})$$

$$\hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MSRes \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}$$

interval for y_0

called as
prediction
interval.

e.g. $x = 40$

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x, n = 120, 17$$

"predict" interval for given $x = 40$ with
95% confidence based on T^* is

$$\left(\hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MSE_{Res} \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}, \right.$$

$$\left. \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MSE_{Res} \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} \right)$$

$$\Rightarrow (83.5, 156.84)$$

CI for mean response $\Rightarrow (109, 132)$
interval length = 23

prediction interval $\Rightarrow (156.84, 83.5)$

interval leng = 73.34

predictⁿ int always wider than CI.
for $x = x_0$ ~~near~~

or mean cannot vary much as compared to value corresp. to given prediction interval.