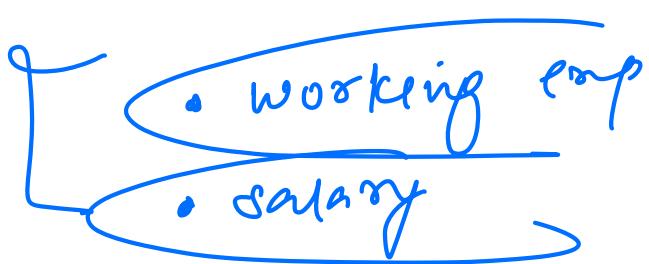


Project 2: Model Stock Prices Using Linear Regression & test all the 5 Assumption of LR Model Using statistical Methods

Linear Regression → statistical model

- linear Regression is a way to find the relationship between 2 or more things.

→ If one thing changes (year of working exp), how does the salary changes.



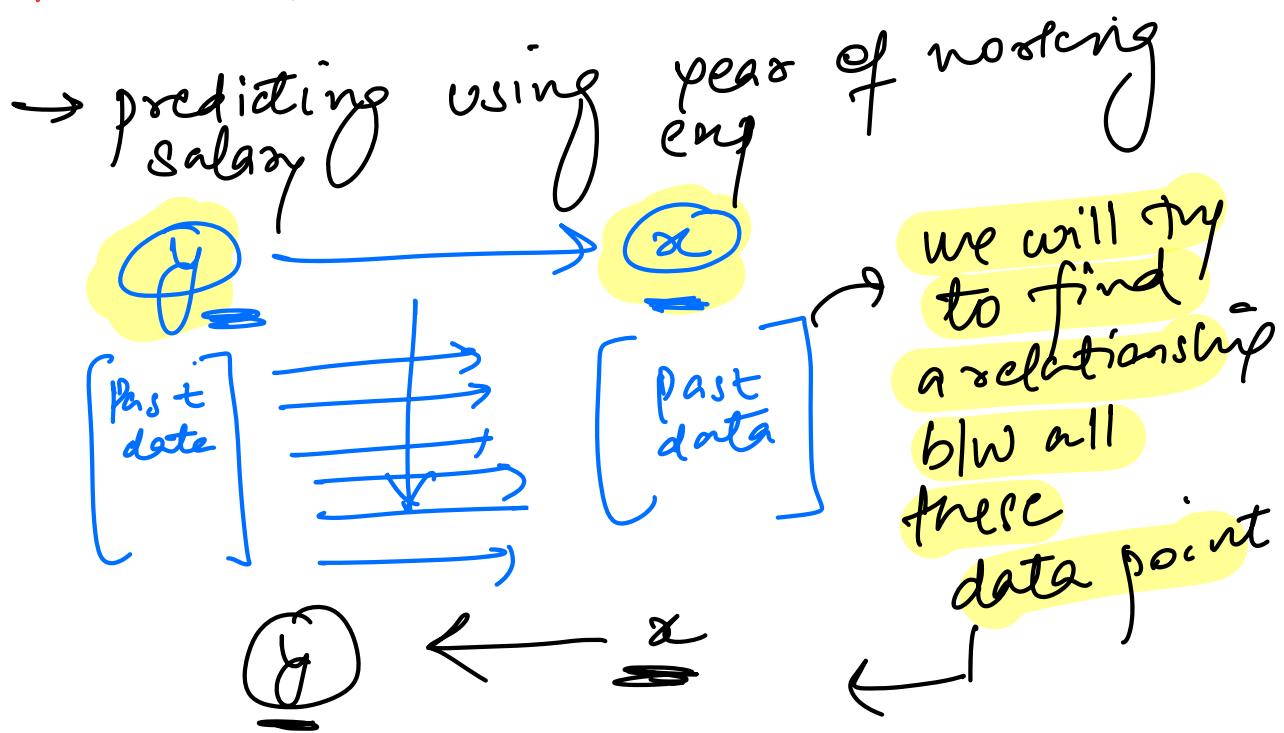
(x) ← relationship b/w these 2 variable..
 (y) ←

[If year of working exp is higher → salary high]

Tone

pattern → Past data that we observed

- LR helps you predict a value based on past data



linear
Regression
Eqⁿ

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

- y = salary
- x = years of working exp
- β_0 = intercept → starting point of y when $x=0$
- β_1 = coefficient of x

NOTE : linear regression tries to draw a "straight line" through the

data point

"Best fitted line"

Goal

→ Predict y

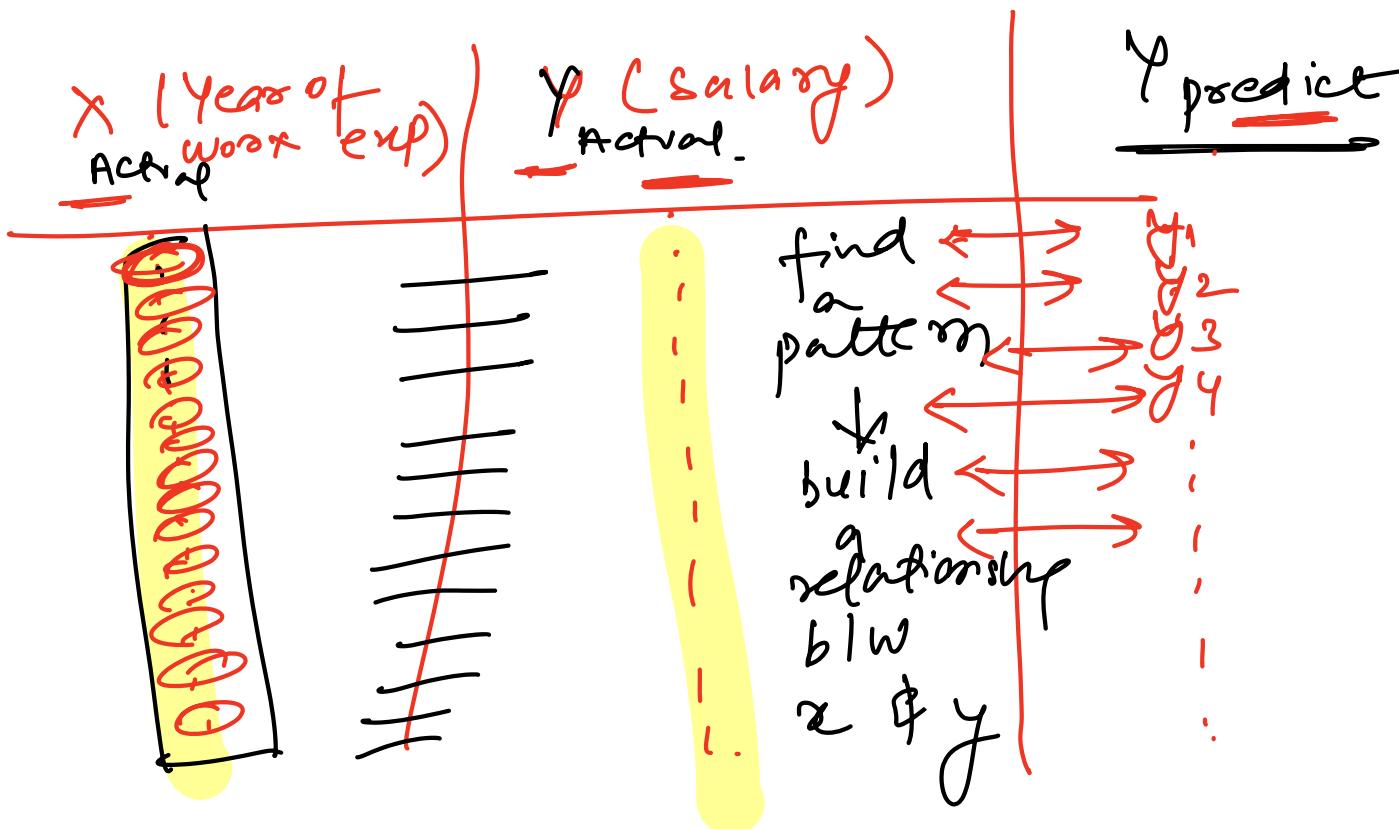
LR find the best fitted line that predict the output as close as possible to the real value.

"Best" Means

→ minimize the error b/w

→ what we predicted

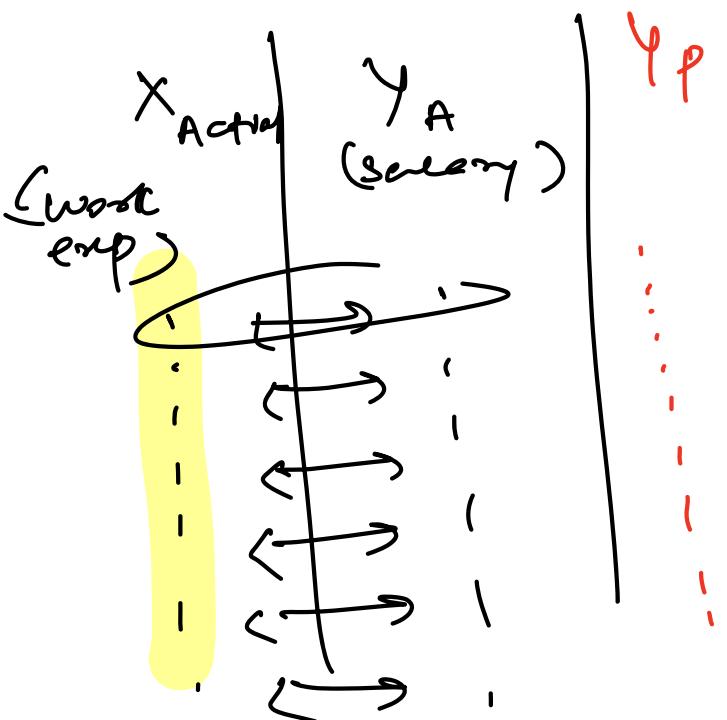
→ what actually happen.



- When we run a LR model
- Train the model

500 working professionals

salary
working exp



How close

Y_A & Y_P are.

$Y_A | Y_P \rightarrow$ close
Best

$\cancel{Y_A} | Y_P \rightarrow$ Bad
far off
 $\$100,000$
 $\cancel{\$50,000}$

[We run a LR model to find a relation b/w x & y]

$$[y = \beta_0 + \beta_1 x] \rightarrow \text{LR Model}$$

Using all these 500 data point

$$\beta_0 = 0.5$$

$$\beta_1 = 2$$

LR Model

$$y = 0.5 + 2x$$

$$\begin{aligned} x = 2 &\rightarrow \begin{array}{|c|} \hline y_1 \\ \hline \end{array} \\ x = 3 &\rightarrow \begin{array}{|c|} \hline y_2 \\ \hline \end{array} \\ x = 3.5 &\rightarrow \begin{array}{|c|} \hline y_3 \\ \hline \end{array} \end{aligned}$$

Linear Regression

→ LR is a statistical method which is used to find a relationship b/w a dependent variable & one or more independent variables.

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

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

y
salary → year of work exp

→ y is dependent on ' x ' was not dependent on anything] Indep
→ x (working exp) →

MBA

depen y → year of exp, technical skills, which university you graduated

1

2

3

x_1

x_2

x_3

$$\text{salary} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

$\beta_0 \rightarrow$ int except

$\beta_1, \beta_2, \beta_3 \rightarrow$ Coeff of x_1, x_2, x_3

$\varepsilon \rightarrow$ error term

↳ captures the variation that is not explained by the model.

~~Type~~

- ✓
- One independent Var (x)

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

↳ Simple

↳ Linear Regression

- More than one Indep. Var (x_1, x_2, \dots)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

$$- y = \beta_0 + \beta_1 x_1 \quad x = y \\ \textcircled{+} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

↳ Multiple Linear Regression

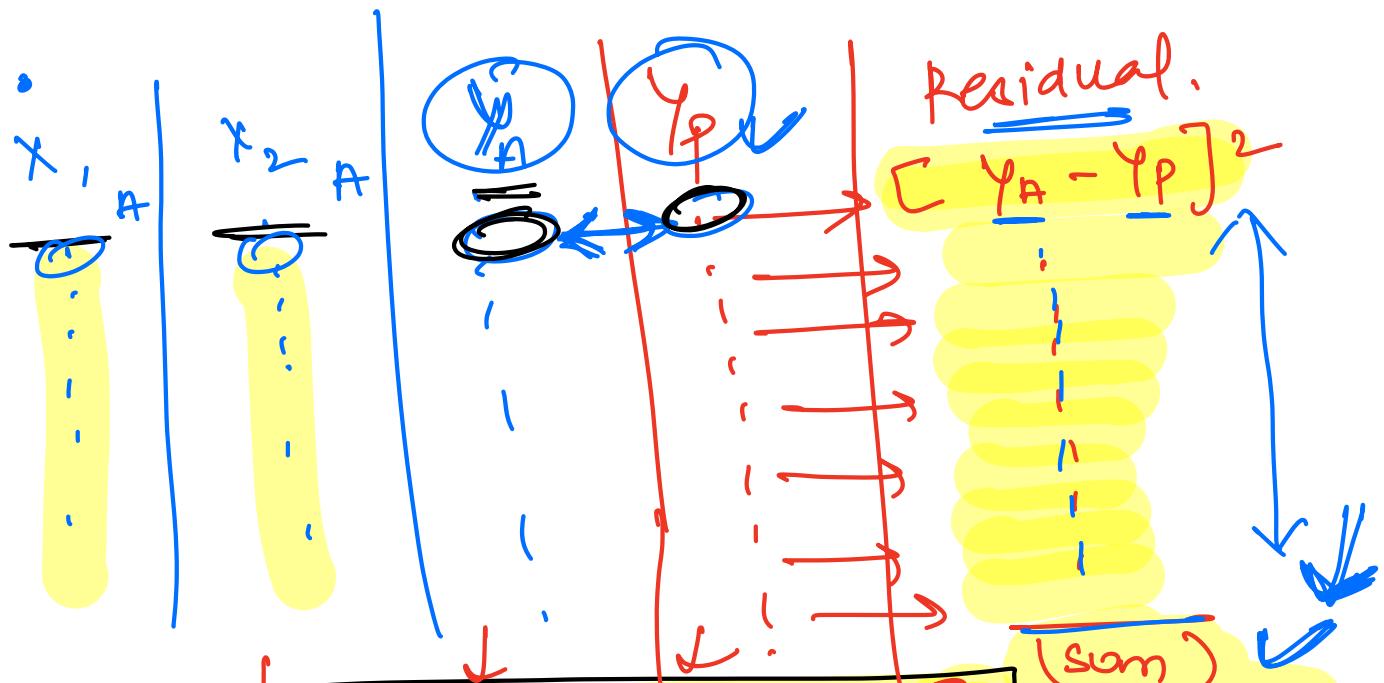
How do you calculate coefficient

$$\underbrace{\beta_0, \beta_1, \beta_2}_{\text{Intercept}} \dots \underbrace{\beta_n}_{\text{Coeff.}}$$

Intercept Coeff.

↳ The way we estimate the parameters is in such a way that the linear regression model minimizes the sum

of square errors [Residual]
 b/w the Actual Value &
 the predicted value.



✓ $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

$\beta_0 = 0.7$

$\beta_1 = 0.5$

$\beta_2 = 1.5$

$\epsilon = 0.05$

↓ Error term

$\hat{Y} = 0.7 + 0.5 X_1 + 1.5 X_2 + 0.05$

final LR model

β_0
 β_1
 β_2
 ϵ

least squares error
 minimize $\rightarrow (y_A - y_P)^2$

→ The technique to find the parameter of LR model is called **least square error method**.

summarise

- statistical model relationship b/w $y \& x(x_1, x_2, \dots)$
- LR tries to find best fitted line.
- Types \rightarrow simple LR (a) $[A \text{ vs } P]$
- Types \rightarrow Multiple LR (x_1, x_2, \dots)
- $B_0, B_1, \dots \rightarrow$ least square error method

====

Assumption of Linear Regression Model

- * * * * *
- what are the Ass. of LR Model
 - how do you test all the Ass. of LR Model.

• LR has 5 assumption

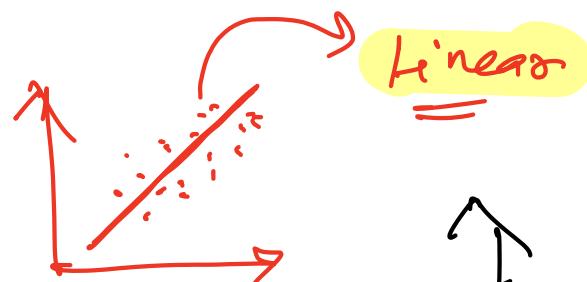
=====

a) **linearity**

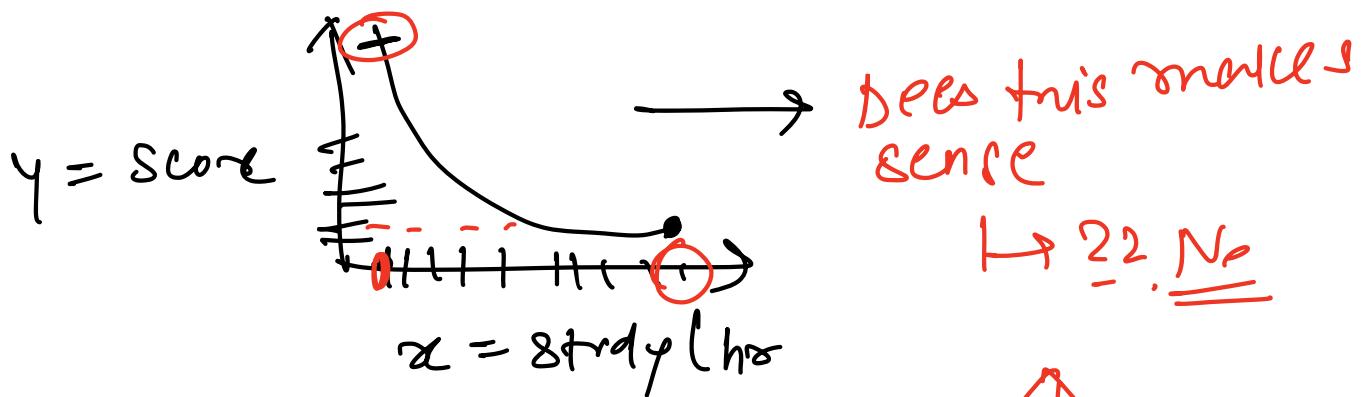
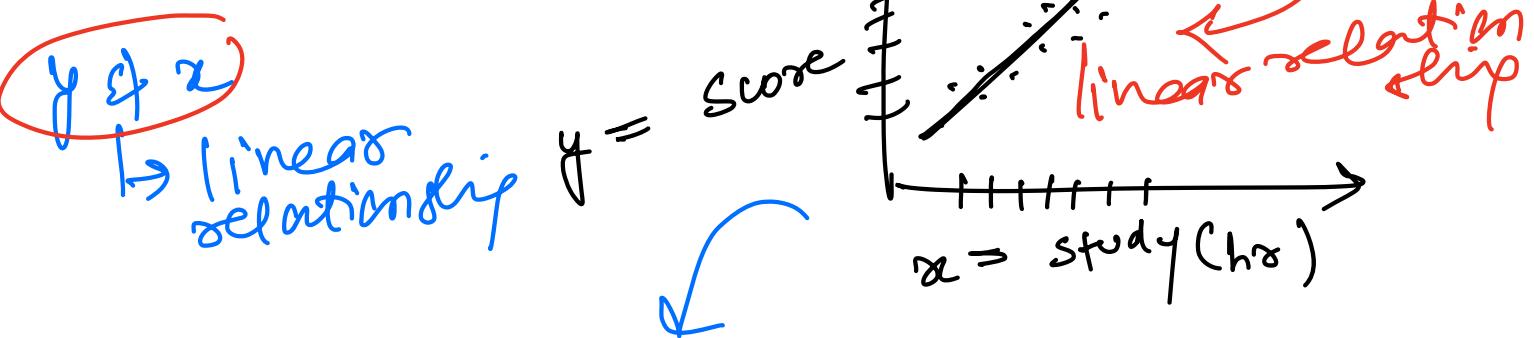
b) Independence of errors

- c) Homoscedasticity
- d) Normal Distribution of Errors
- e) No Multi⁹ collinearity

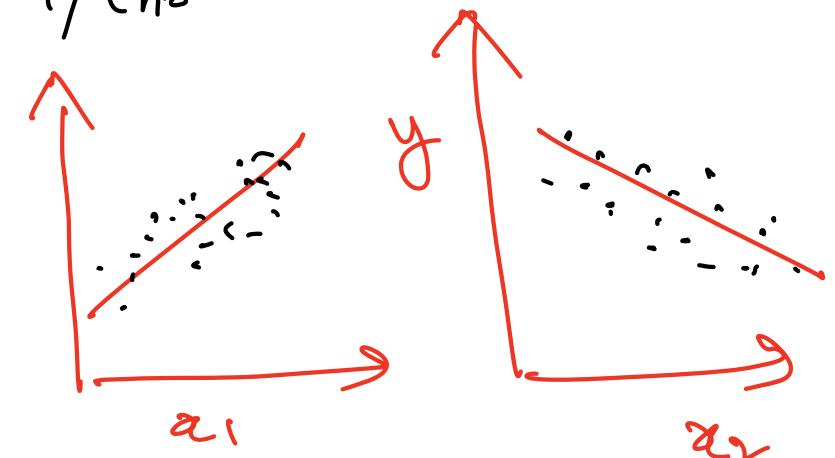
(a) Linearity: There should be a linear relationship b/w your dependent & independent variable



straight line



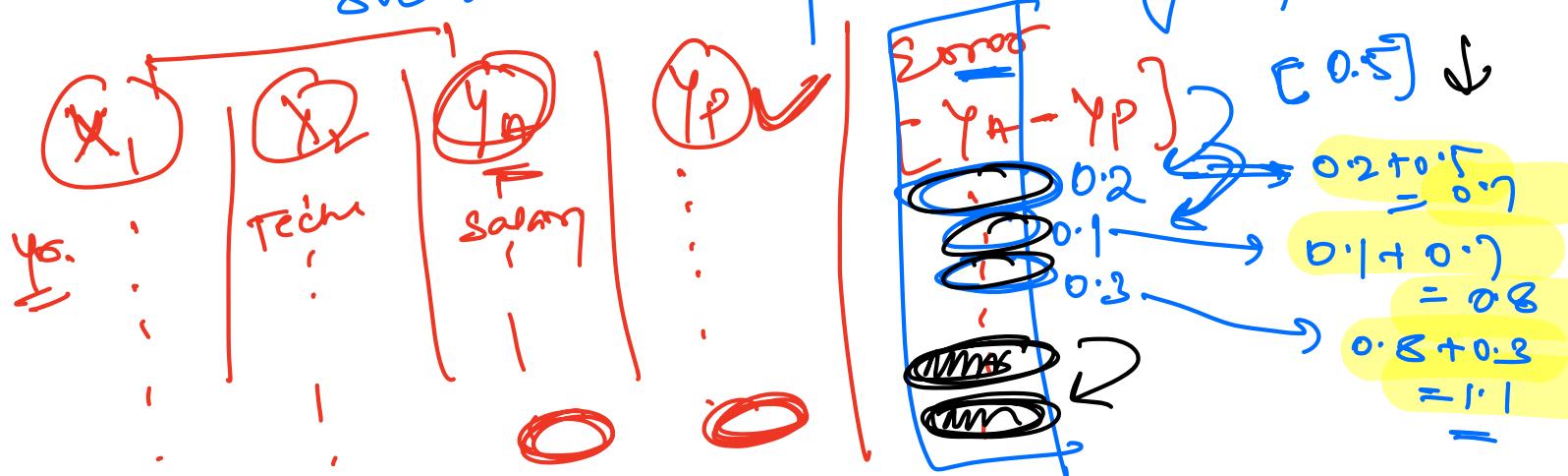
y [x_1
 x_2]
→ Data point



Q2 Independence of Errors (Residual)

- The residual errors are independent of each other.

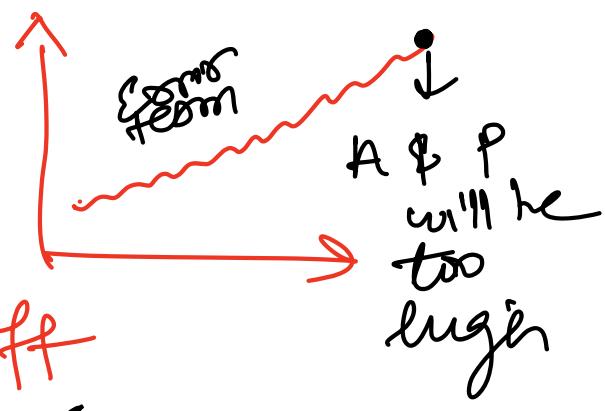
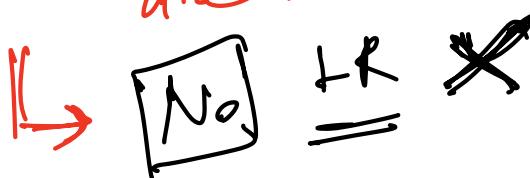
Mean: The error (difference b/w A & P) should not follow any pattern



- If it's dependent

If error is high

$\rightarrow A \& P$ are too off



- All the error (Diff A vs P) should be independent of each other

Eg: By mistake one person's salary prediction is too high \Rightarrow Diff A & P

This should not affect
the prediction of next
person's salary.

error higher

③ Homoskedasticity (Constant Variance)

The variance of residual is
constant across all levels of independent
variable.

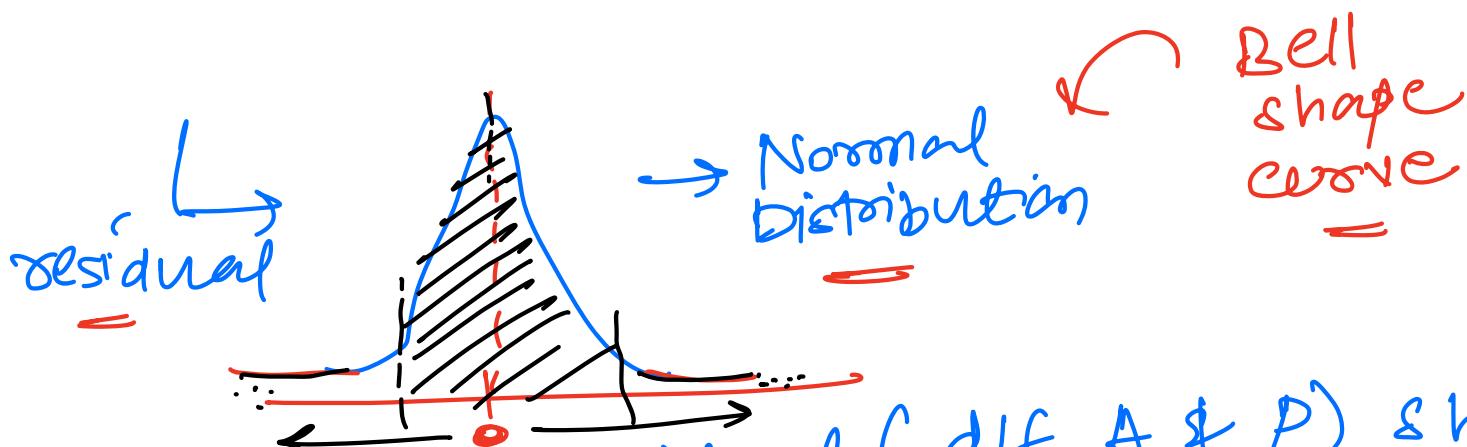
Means : The spread of the error should
be roughly the same across
all values of IP.

Eg : When you are predicting the
salary for someone who has
140 of work exp & someone
with 200 of work exp , the
error in prediction should
be similar

- small for one
- Big for the other

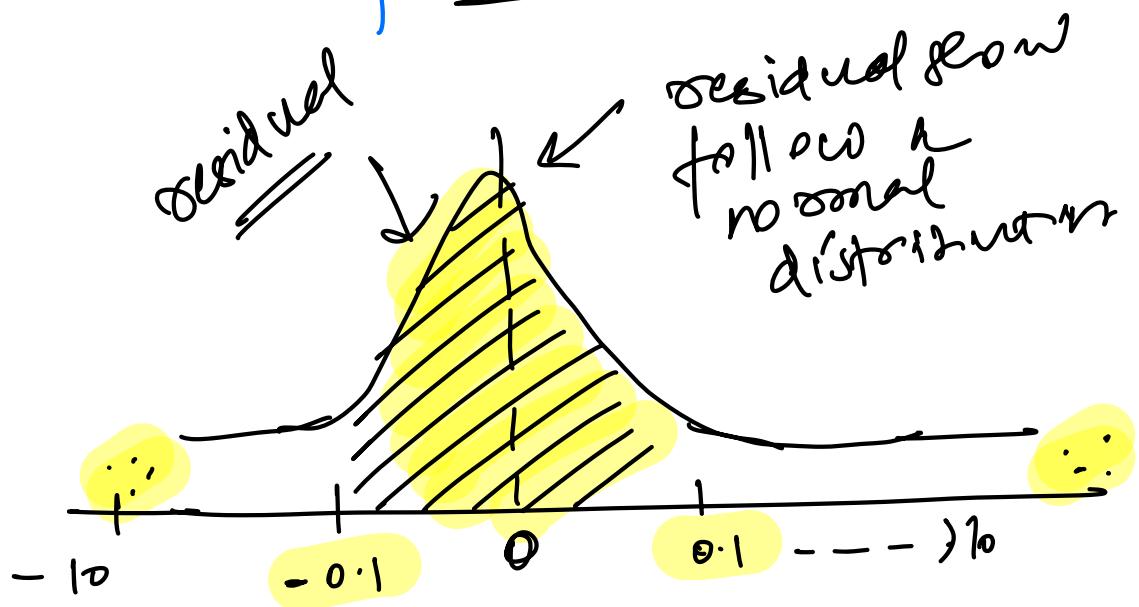
④ Normal Distribution of Errors

The residual should be normally distributed



Means: The residual (diff A & P) should follow a bell shape curve.

→ most errors should be small
few errors can be large.

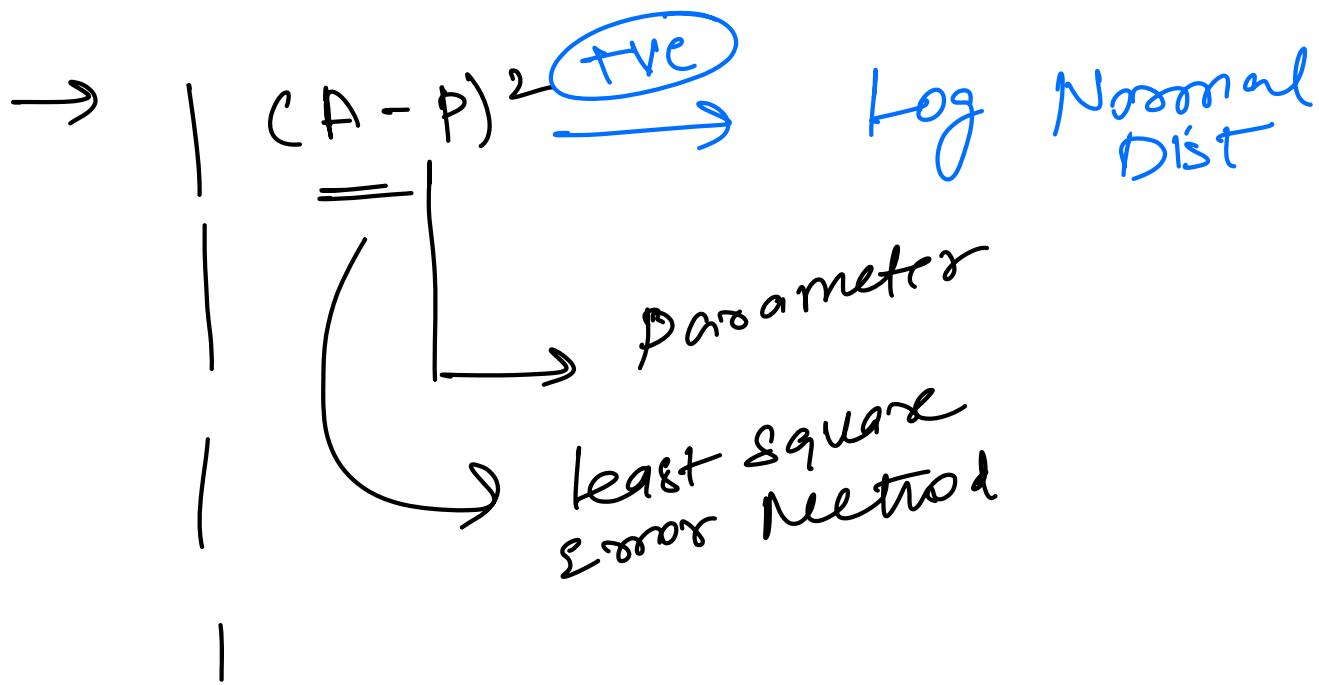
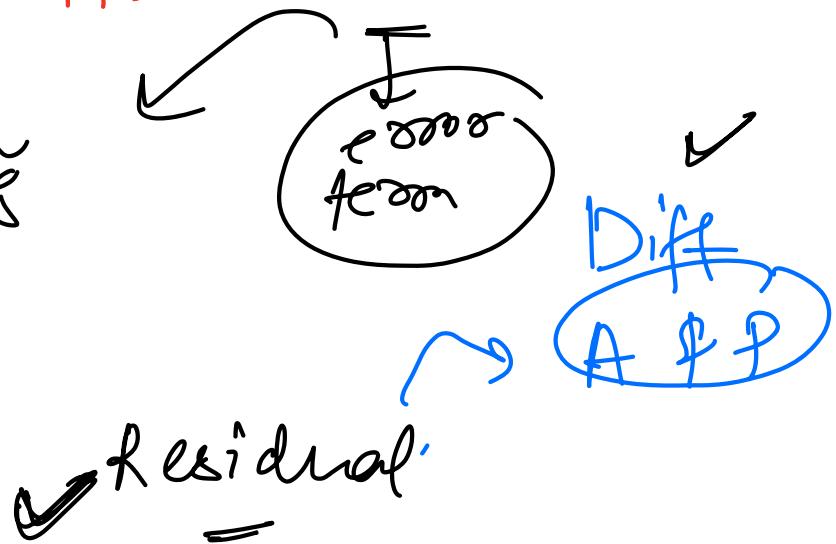


Error Term

} Residual

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

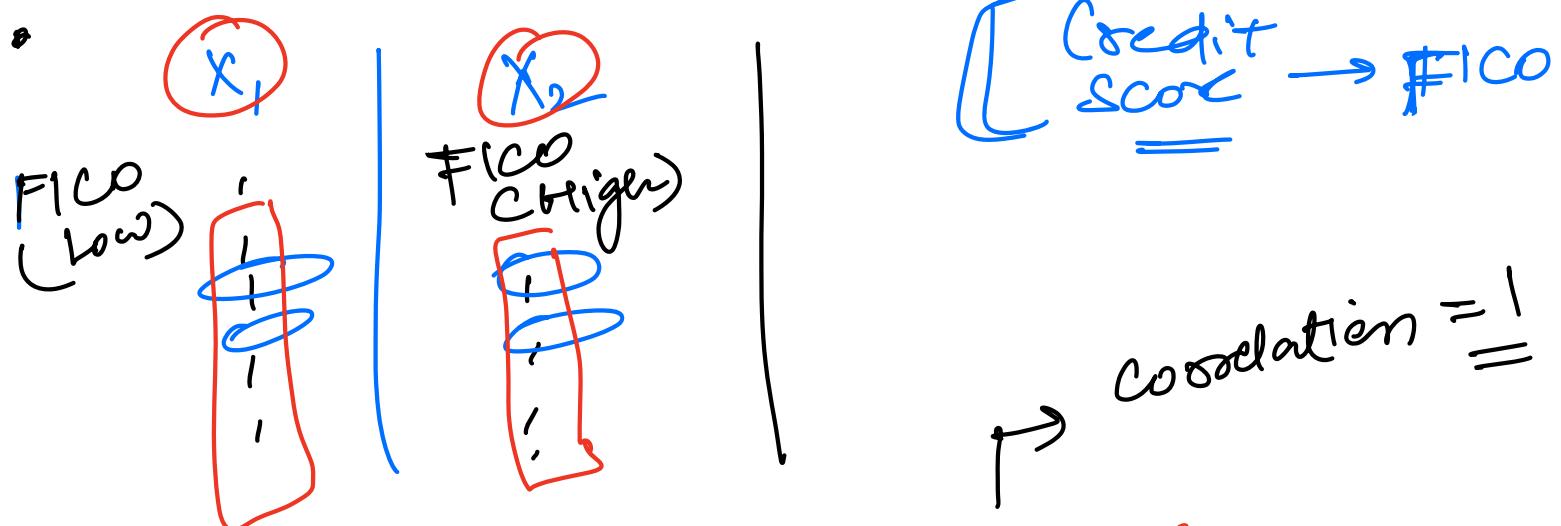
✓ error which the model is not able to explain



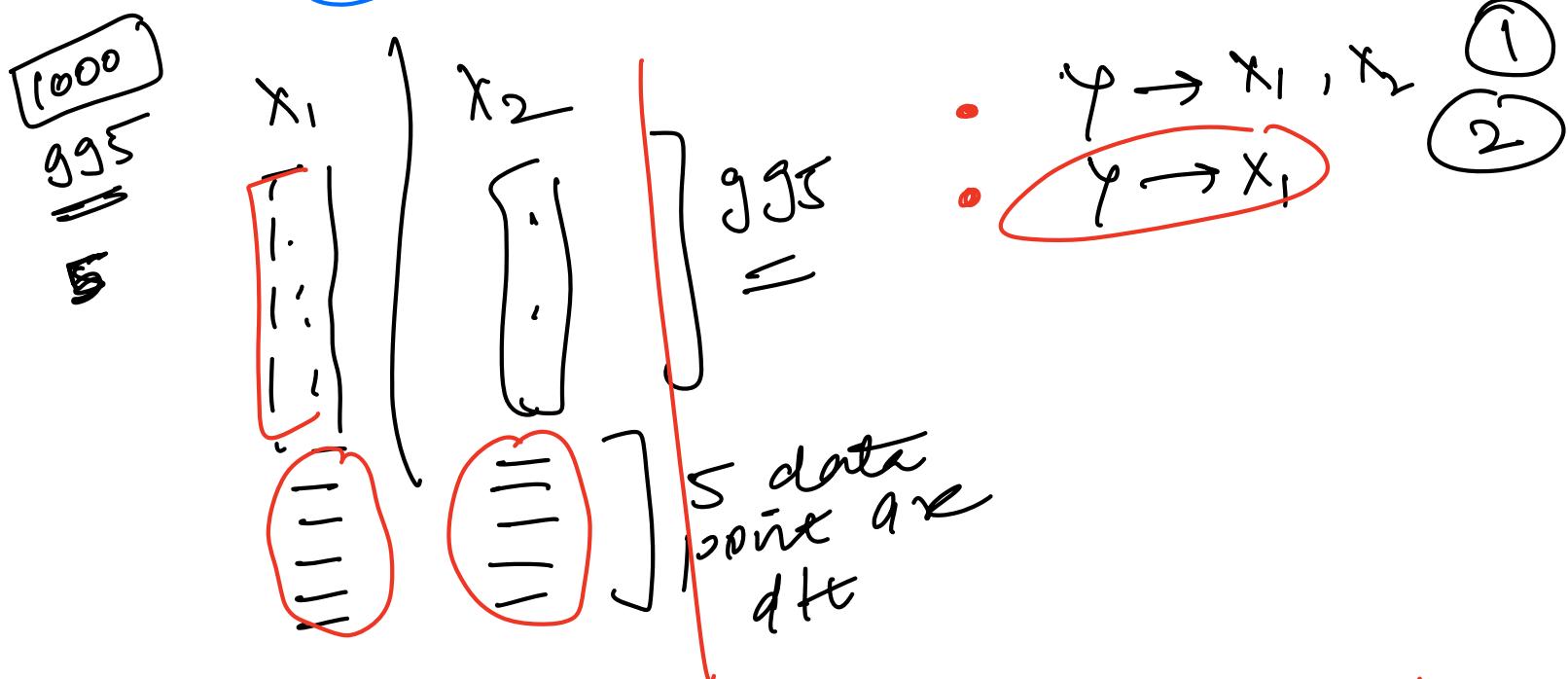
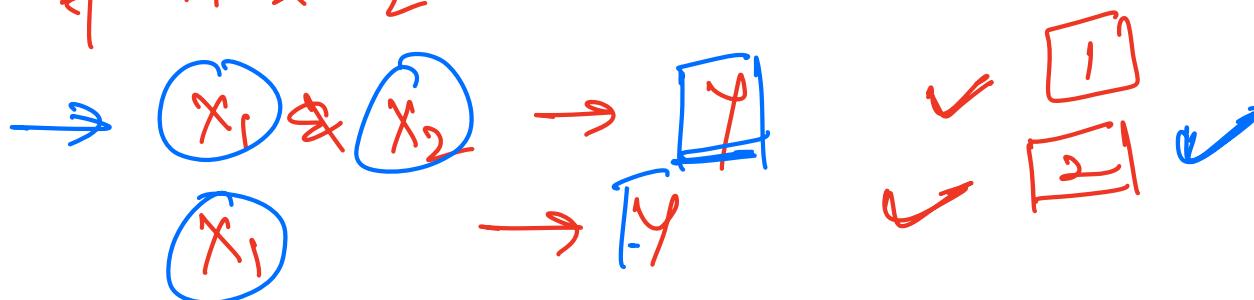
plot \rightarrow Residual $\Rightarrow (A - P) \rightarrow$ Normal dist'n

5 No Multicollinearity

- The independent vars ($x_1, x_2, x_3 \dots$) should not be highly correlated with each other

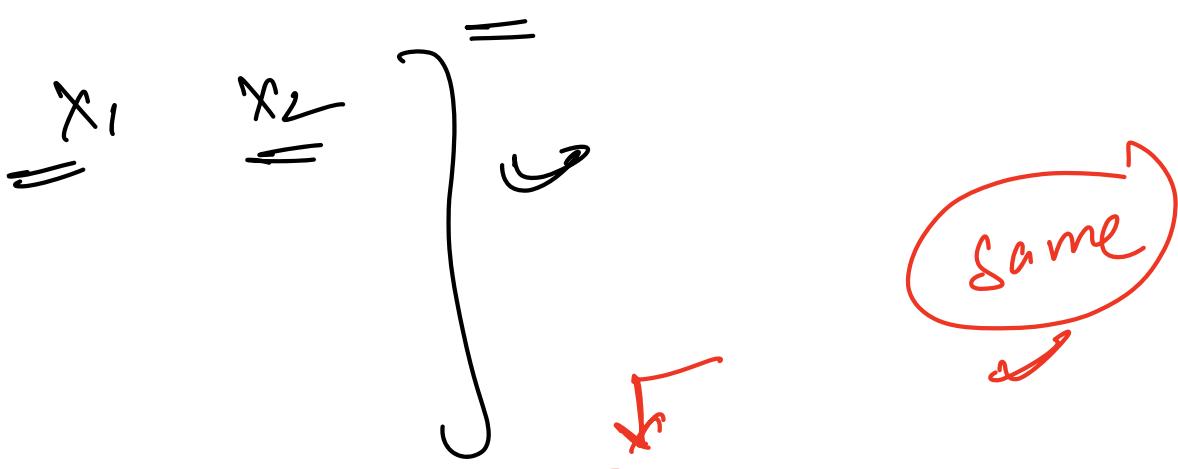


- If x_1 & x_2 are exactly same



- When if vars are highly correlated
→ it makes no sense in having those vars in our LR model
- We need to make sure if variable should not be closely related =

to each other



House
price

House
size

size of each
room + size
of all remaining
space

If the i/p vars are closely aligned
with each other \rightarrow it makes no sense
of having both of
them together

\hookrightarrow Drop 1 var

SIR \rightarrow α
MLR \rightarrow α_1, α_2

Summarize

① Linearity : Linear relationship b/w dep & indep var
(y) (x_1, x_2, \dots).

② Independence of Errors : The residual should be independent of each other

③ Homoscedasticity : The variance of residual should be const

④ Normal Distribution of Errors → The residual should be normally distributed

⑤ No multicollinearity → The independent var (x_1, x_2, \dots) should not be highly correlated with each other.

→ → → → → → → → → →