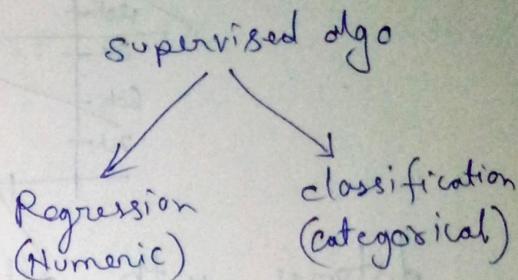
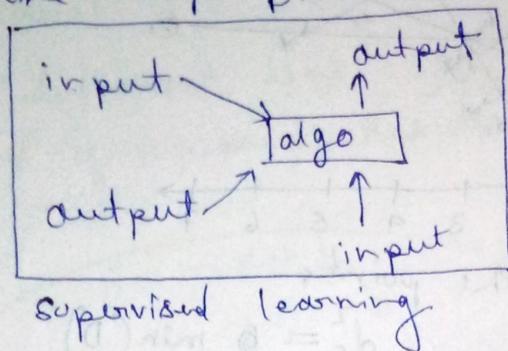


01 Oct 2023

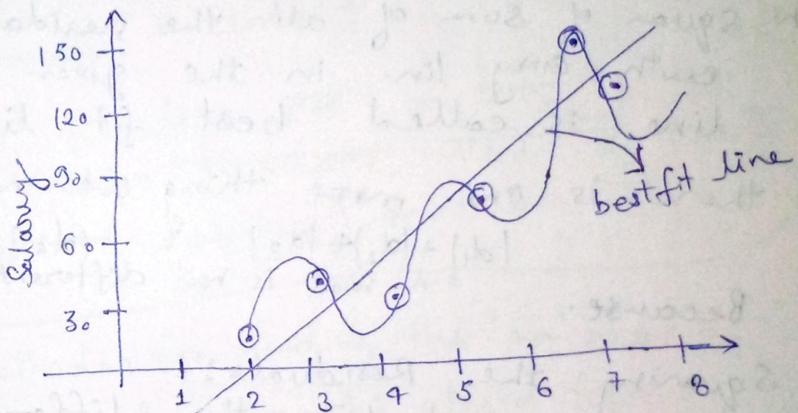
# Linear Regression

It is algorithm that is used to solve supervised machine learning problem.



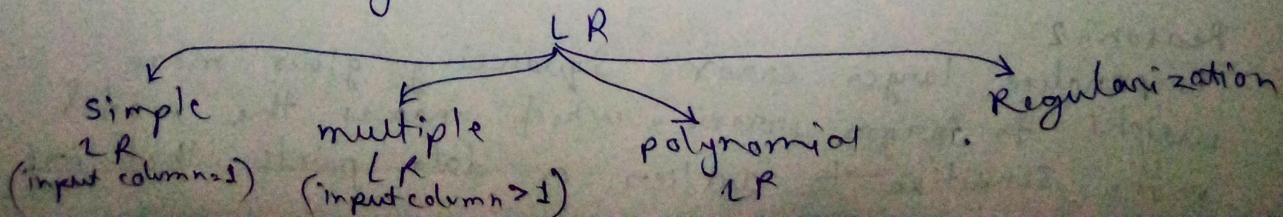
experience	Salary
5.6	75000
7.2	120000
3.1	47000
4.4	37000
2.0	16000
:	
7	140000

training data



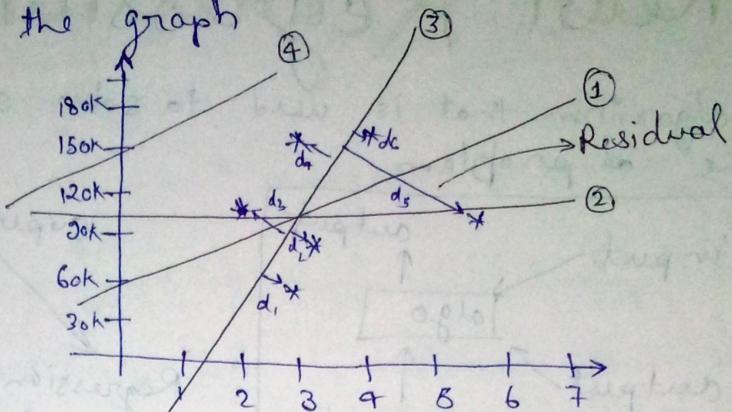
- \* machine learning Algorithm को use करके हम इसे model तैयार करते हैं माले ही वह Training set data का accurate work ना कर सकता तो किन वह Testing data पर accurately work करता चाहिए

- \* The user may be perfectly work with training data but it not perfectly work with testing data hence our end goal is to make the model which is perfectly work with testing data provided by user.



How to select bestfit line?

The best fit line is always most closer to each point on the graph



distance with all the points

$$d_1^2 + d_2^2 + d_3^2 + \dots + d_6^2 = \min(D)$$

\*the square of sum of all the residual must be minimum with any line in the given then that ~~that~~ line is called best fit line

\* there is one more thing we not use mod  $|d_i|$

$$(d_1) + (d_2) + (d_3) \dots + (d_6) = \min(D)$$

Because mod is not differentiable at zero

Squaring the Residuals:

when we calculate the difference between the observed Y-values and the predicted Y-values (residuals), we square these differences for this reason.

Reason - 1

To Remove the sign: Squaring ensures that all residuals are positive. This is important because if we only added up the differences without squaring them, the positive and negative residuals would cancel each other out, and we wouldn't be able to assess the overall goodness of fit.

Reason - 2

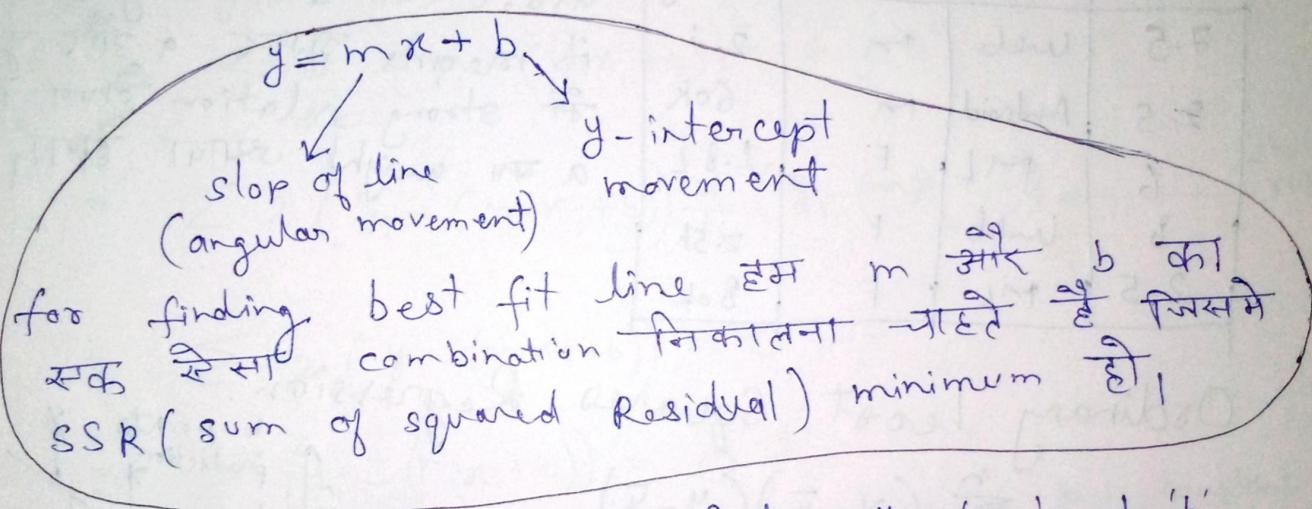
To penalize larger error: Squaring gives more weight to larger errors, which makes the model more sensitive to outlier or data points that are

that have a significant impact on the overall fit of the line.

$$\boxed{d_1^2 + d_2^2 + d_3^2 \dots d_n^2 = D}$$

In the above  $D$  is called sum of squared Residual.

for different line  $SSR_1, SSR_2, SSR_3, \dots, SSR_n$   
(sum of squared residual)



there are two methods for finding the 'm' and 'b'

(i) OLS (ordinary least square): closed form  
it is set of formula (statistical method)

(ii) Gradient Descent (calculus method):  
Non closed form. also called approximation  
based technique., it is used for higher  
dimension dataset. SGD Regression class.  
Home we use.

## Regression Intuition:

for a 2D line:  $y = mx + b$

for a 3D line:  $y = ax_1 + bx_2 + c$

for a 4D line  $y = ax_1 + bx_2 + cx_3 + d$  offset

$x_1$	$x_2$	$x_3$	$y$
Exp	profile	Gender	Salary
7.5	Web	M	2.1
2.5	Android	M	60k
6	ML	F	1.8 L
8	Web	F	2.5k
2.5	ML	F	80k

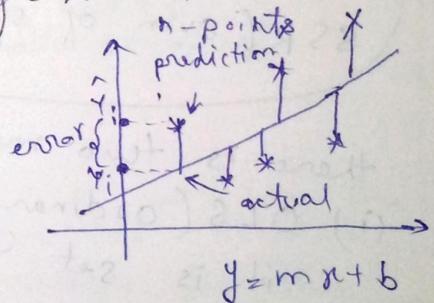
\* in the above,  $a, b, c$  is weight  
 \* if we calculate  $y$  from  $a, b, c$  then  $a, b, c$  are weight  
 it means there is a strong relation between  $x_1, x_2, x_3$  and  $y$   
 \* strong relation द्वारा तब a का weight ज्यादा होता।

## Ordinary Least Squares Regression.

$$m = \frac{\sum_{i=0}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=0}^n (x_i - \bar{x})^2}$$

where  $\bar{x}, \bar{y}$  = mean value.

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



~~proof:~~

$$\text{sum of square} = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_n - y_n)^2$$

$$E = SS = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \text{error function}$$

$$\therefore \hat{y}_i = mx_i + b$$

$$E(m, b) = SS = \sum_{i=1}^n (y_i - (mx_i + b))^2$$

the error function depends on the two values  $m, b$

in the above equation. we have to calculate partial differentiation w.r.t. m and b minimization.

$$\frac{\partial}{\partial b} \sum (y_i - (mx_i + b))^2 = 0 \quad \text{--- (i)}$$

$$\frac{\partial}{\partial m} \sum (y_i - (mx_i + b))^2 = 0 \quad \text{--- (ii)}$$

for equation (i)

$$\frac{\partial F}{\partial b} = \sum \frac{\partial}{\partial b} (y_i - (mx_i + b))^2 = 0$$

$$= \sum 2(y_i - (mx_i + b)) \frac{\partial}{\partial b} (y_i - (mx_i + b)) \quad \text{chain rule}$$

$$= \sum -2(y_i - (mx_i + b)) = 0$$

$$= \sum (y_i - (mx_i + b)) = 0$$

$$= \sum (y_i - mx_i - b) = 0$$

$$= \sum y_i - \sum mx_i - \sum b = 0$$

$$= \sum y_i - m \sum x_i - nb = 0$$

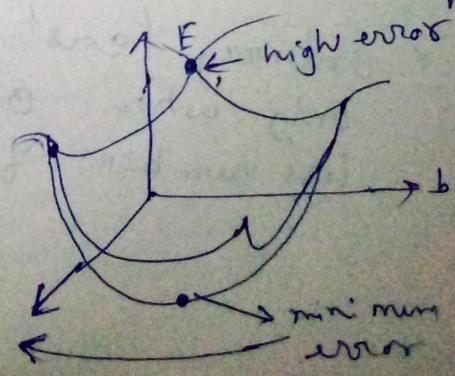
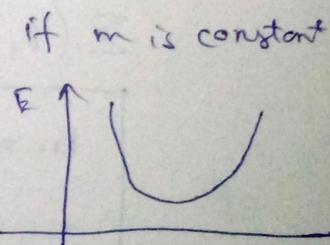
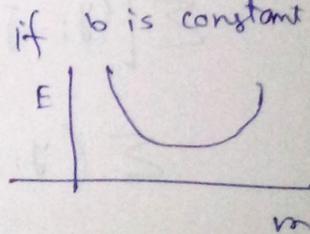
$$= \frac{\sum y_i}{n} - \frac{m}{n} \sum x_i - b = 0$$

$$= \bar{y} - m \bar{x} - b = 0$$

$$= \boxed{b = \bar{y} - m \bar{x}}$$

$$\text{or } b = \frac{\sum y - n \bar{x}}{n}$$

for calculating slope = zero  
derivative = zero



for equation (ii)

$$\frac{\partial F}{\partial m} = \frac{d}{dm} \sum (y_i - (mx_i + b))^2 = 0$$

$$\sum \frac{d}{dm} (y_i - (mx_i + \bar{y} - m\bar{x}))^2 = 0$$

$$\sum 2(y_i - (mx_i + \bar{y} - m\bar{x})) \cancel{\cdot} \cancel{m}$$

$$\sum 2(y_i - (mx_i + \bar{y} - m\bar{x})) - (x_i - \bar{x}) = 0 \quad \text{by chain rule}$$

$$-2 \sum (y_i - (mx_i + \bar{y} - m\bar{x})) (x_i - \bar{x}) = 0$$

$$\sum (y_i - (mx_i + \bar{y} - m\bar{x})) (x_i - \bar{x}) = 0$$

$$\sum [(y_i - \bar{y}) - m(x_i - \bar{x})] (x_i - \bar{x}) = 0$$

$$\sum [(y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2] = 0$$

$$\sum (y_i - \bar{y})(x_i - \bar{x}) - m \sum (x_i - \bar{x})^2 = 0$$

$$\sum (y_i - \bar{y})(x_i - \bar{x}) = m \sum (x_i - \bar{x})^2$$

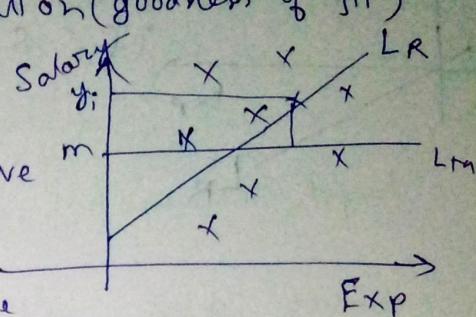
$$m = \left[ \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \right] \quad \text{or} \quad \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

- \* Ordinary least squares Regression is used only when our data set is small and less number of columns.

## Regression Matrices:

- (i)  $R^2$  - coefficient of Determination (goodness of fit)
- (ii) Adjusted  $R^2$

\*  $L_M$  line is used when we have  
Not a experience column.



\*  $L_R$  line is used when we have  
two columns Salary and Experience.

\*  $R^2$  is used to find How much better  $L_R$  from  
 $L_M$

$$R^2 = \frac{\text{Var}_m - \text{Var}_R}{\text{Var}_m}$$

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

$$\text{Var}_m = SSR_m = \sum_{i=1}^n (y_i - \hat{y}_m)^2$$

$$\text{Var}_R = SSR_R = \sum_{i=1}^n (y_i - \bar{y})^2$$

in the above.  $0 < R^2 < 1$

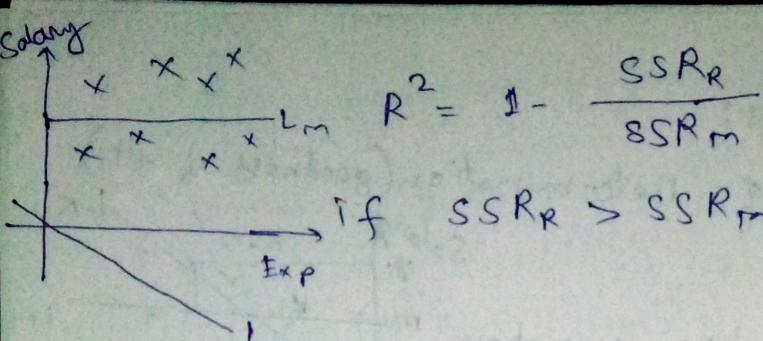
$R^2$  is calculated in % example 60%.

80% of  $R^2$  means. salary and experience  $\rightarrow$  correlation ~~are~~ used ~~that~~ we are only able to explain the 80% variation in data.

Number of column  $\propto R^2$

interview question - What if  $R^2$  -ve & how is it?

Answer - Yes



if  $SSR_R > SSR_m$ ,  $R^2$  is -ve

$$R^2_{adj} = 1 - \left[ \frac{(1-R^2)(n-1)}{(n-k-1)} \right]$$

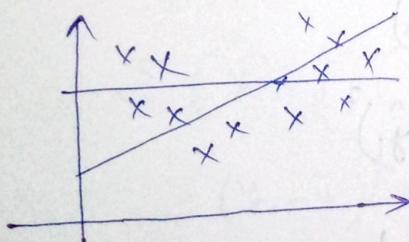
$n$  = Number of rows in data.

$k$  = Number of independent columns

### Polynomial Regression:

$$y = \beta_0 + \beta_1 x \xrightarrow{\text{simple linear Regression.}}$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n \xrightarrow{\text{multiple linear Regression}}$$



x	y
35	100

degree 2  
 $x^0 = 1, x^1 = 35, x^2 = 1225, \dots$

Hence

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

- \* if degree is very high then overfitting
- \* if degree is very low then underfitting

$x_1$	$x_2$	$y$

input column is true for degree = 2

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2$$

### 5. Adjusted R<sub>2</sub> score:

$$R^2_{\text{adjusted}} = 1 - \left[ \frac{(1 - R^2)(n-1)}{(n-1-k)} \right]$$

n = number of rows.

k = total number of independent input column

- \* if irrelevant R<sub>2</sub> score is added then Adj R<sub>2</sub> score is decrease
- \* if relevant column is added then adj. R<sub>2</sub> score is increase.
- \* adjusted R<sub>2</sub> score is useful when we have a multiple linear regression.

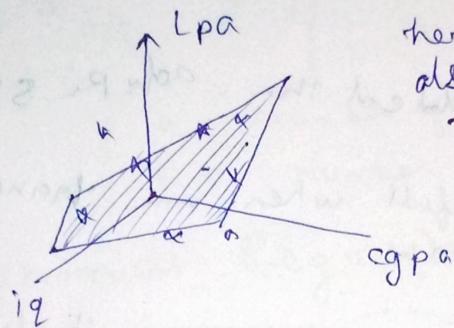
interview question

Ques: between R<sup>2</sup> and adj R<sup>2</sup> which will always be bigger

# Multiple Linear Regression

- it have more than one independent columns.
- example

cgpa	gender	iq	package (LPA)
8	Male	124	3.8
7		120	3.5
6		115	3.2
5		110	2.8



here we draw the plane.  
also called best fit plane  
this is for 3-D.  
for the 4-D hyperplane.

for 2D Data :  $y = mx + b$

for 3D Data :  $y = mX_1 + mX_2 + b$   
OR

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

for 4D Data :  $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$

for nD Data :  $y = \beta_0 + \sum_{i=1}^n \beta_i X_i$

here coefficient  $\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are weights  
e.g.  $y$  calculate ~~and~~  $\beta_1 X_1$  has weight  $\beta_1$

$$lpa = \beta_0 + \beta_1 \text{cgpa} + \beta_2 \times iq$$

$$\text{if } \beta_1 > \beta_2$$

cgpa has more weight than iq

for the calculating lpa

$\beta_0$  = offset / intercept value

predicted	c gpa	iq	gender	ipa	actual
	$x_{11}$	$x_{12}$	$x_{13}$	$y$	
	$x_{11}$	$x_{12}$	$x_{13}$		

$$\rightarrow \hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

for calculating  $\hat{y}$  (predicted)

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 x_{11} & \beta_2 x_{12} & \beta_3 x_{13} \\ \beta_0 & \beta_1 x_{21} & \beta_2 x_{22} & \beta_3 x_{23} \\ \vdots & \vdots & \vdots & \vdots \\ \beta_0 & \beta_1 x_{m1} & \beta_2 x_{m2} & \beta_3 x_{m3} \end{bmatrix} = \hat{Y}$$

if we have n rows and m columns

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 x_{11} & \beta_2 x_{12} & \cdots & \beta_m x_{1m} \\ \beta_0 & \beta_1 x_{21} & \beta_2 x_{22} & \cdots & \beta_m x_{2m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \beta_0 & \beta_1 x_{n1} & \beta_2 x_{n2} & \cdots & \beta_m x_{nm} \end{bmatrix}$$

$$\hat{Y} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \cdots & x_{1m} \\ 1 & x_{21} & x_{22} & x_{23} & \cdots & x_{2m} \\ 1 & 1 & 1 & 1 & \ddots & \vdots \\ 1 & 1 & 1 & 1 & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \cdots & x_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix}$$

$$\hat{Y} = X\beta$$

$X$  = prediction of all input values matrix

$\beta$  = all the coefficient matrix

$\hat{Y}$  = prediction matrix

$$Y(\text{actual}) = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad e = Y - \hat{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

in simple linear regression

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

$$E = e^T e$$

$$E = [(y_1 - \hat{y}_1) \ (y_2 - \hat{y}_2) \ \dots \ (y_n - \hat{y}_n)]_{1 \times n} \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}_{n \times 1}$$

$$\begin{aligned} E &= (Y - \hat{Y})^T (Y - \hat{Y}) \\ &= (Y^T - \hat{Y}^T) (Y - \hat{Y}) \quad \because (A - B)^T = A^T - B^T \\ &= [Y^T - (X\beta)^T] (Y - X\beta) \\ &= Y^T Y - Y^T X \beta - (X\beta)^T Y + (X\beta)^T (X\beta) \end{aligned}$$

proof  $Y^T X \beta = (X\beta)^T Y$

$$\text{let } Y = A \quad X\beta = B$$

$$A^T B = B^T A \quad \text{(i)}$$

$$\because (A^T B)^T = B^T A \quad \text{(ii)}$$

$$A^T B = (A^T \beta)^T$$

$$\text{let } A^T B = C \quad \text{prove } C = C^T$$

$$\text{or } (\gamma^T \times \beta)^T = \gamma^T \times \beta$$

$$E = \gamma^T \gamma - 2 \gamma^T \times \beta + \beta^T \times \gamma^T \times \beta$$

$$\frac{\partial E}{\partial \beta} \rightarrow \frac{\partial}{\partial \beta} [\gamma^T \gamma - 2 \gamma^T \times \beta + \beta^T \times \gamma^T \times \beta]$$

$$\frac{\partial E}{\partial \beta} = 0 - 2 \gamma^T \times \beta + \frac{\partial}{\partial \beta} [\beta^T \times \gamma^T \times \beta] = 0$$

$$\therefore y = A^T A$$

$$\frac{dy}{dA} = 2A^T$$

$$= -2 \gamma^T \times \beta + 2 \gamma^T \times \beta^T = 0$$

$$\cancel{\gamma^T \times \beta^T} = \cancel{\gamma^T \times}$$

$$\beta^T = \gamma^T \times (X^T X)^{-1}$$

$$(\beta^T)^T = [\gamma^T \times (X^T X)^{-1}]^T$$

$$\beta = [(X^T X)^{-1}]^T (\gamma^T \times)^T$$

$$\beta = [(X^T X)^{-1}]^T \cdot X^T \gamma$$

$$\boxed{\beta = (X^T X)^{-1} X^T \gamma}$$

$\therefore X$  is square matrix

$$X = X_{\text{train}}$$

$$Y = Y_{\text{train}}$$

for calculating coefficient and offset we use gradient descent rather than OLS because computational complexity for calculating  $X^T$  or  $\gamma^T$  is very high.

in multiple linear regression, the random errors ( $e$ ) are assumed to be

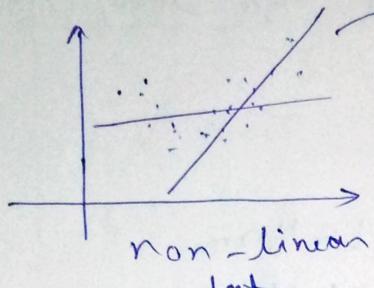
- independently and identically distributed
- normally distributed with zero mean
- constant variance ( $\sigma^2$ )

Date - 17-02-24

# polynomial <sup>linear</sup> Regression.

$y = mx + c$  or  $y = \beta_0 + \beta_1 x \rightarrow$  simple linear Regression.

$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n \rightarrow$  multiple linear Regression



Here linear regression also work but result are not good. Hence here we use polynomial linear regression.

Hence for solve this we add polynomial term in the given data.

example

X	Y	$x^0$	$x^1$	$x^2$	Y
35	100	1	35	1225	100

degree = 2  $\rightarrow$  hyperparameters

For polynomial Regression

for degree = 2 :  $y = \beta_0 + \beta_1 x + \beta_2 x^2$

if column  $\begin{bmatrix} x_1 \\ x_2 \\ Y \end{bmatrix}$  :  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2$

Here  $y \propto \beta_0, \beta_1, \beta_2, \beta_3, \beta_4$   
linearly dependent

Here we called it as polynomial linear Regression

