

SIMPLE LINEAR REGRESSION (supervised model)

> one independent feature and one dependent feature.

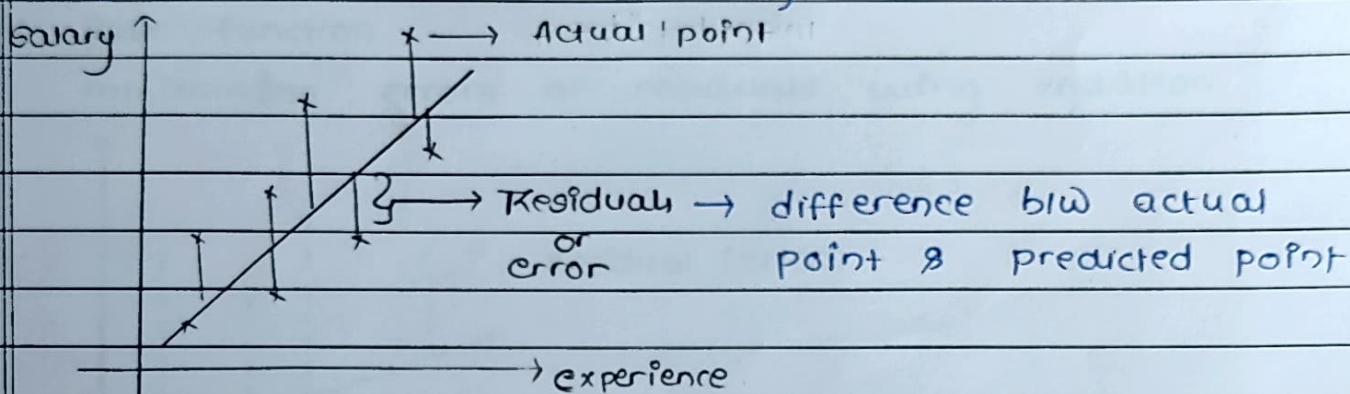
> Eg: Aim to create a model which take the input as height and predict weight.

Variable: height, weight

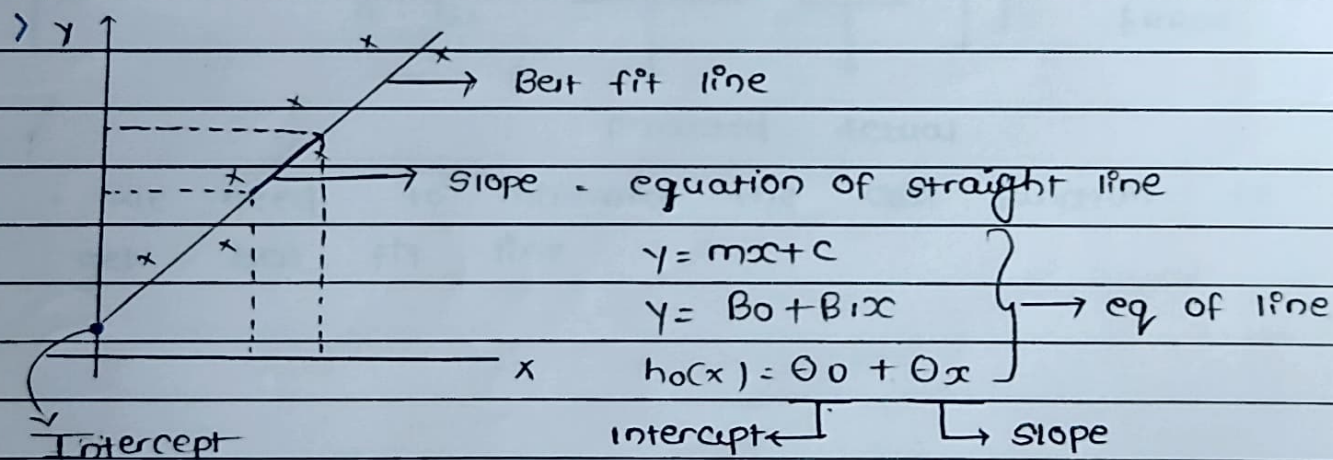
> eg: Aim to create a model which takes no. of rooms and predict the price.

Variable: no. of rooms, price

> Eg: Aim to predict salary based on experience.



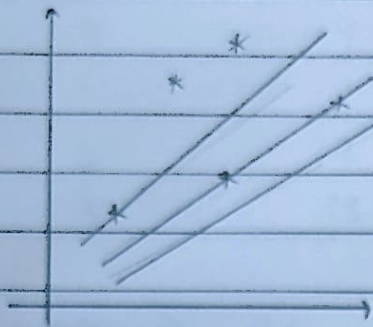
> Based on the training dataset, it find the best fit line in such way that the sum of difference between actual point & predicted point should be minimum



> Intercept: When $x=0$ line meeting the y -axis and that point is known as intercept.

> slope: with the unit movement in the x -axis what is the movement in y -axis

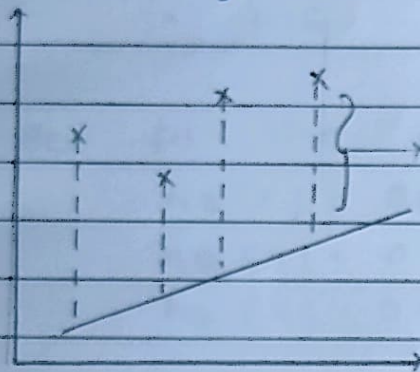
> In order to find best fit line we need to change the values θ_0 and θ_1 .



→ changing the values of θ_0, θ_1 we will get best fit line

> Cost function

minimizing errors or residuals using equation.



Residual / errors

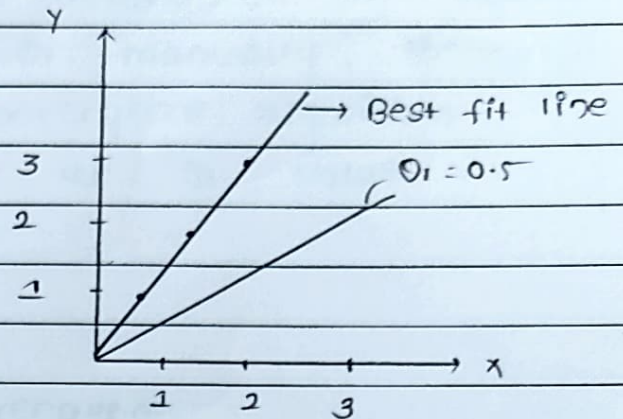
<p>cost function</p> $J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n \left(\underbrace{h(\theta_0, \theta_1)^{(i)}}_{\text{predicted}} - \underbrace{y^{(i)}}_{\text{Actual}} \right)^2$	<p>} → MEAN SQUARE ERROR</p>
---	------------------------------

• We need to minimize the cost function to get best fit line

> eg

Training dataset

x	y
1	1
2	2
3	3



when $\theta_1 = 0.5$, $\theta_0 = 0$

$$J(\theta_1) = \frac{1}{3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

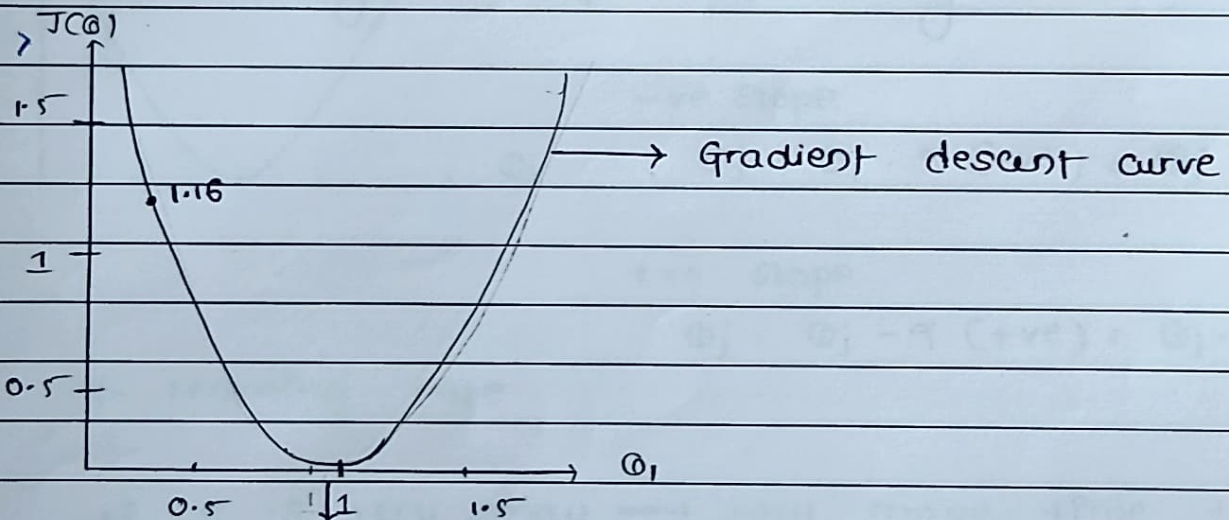
$$= 1.16$$

> Best fit line, $\theta_1 = 0.5$

$$h_0(x) = 0 + 0.5(1) = 0.5$$

$$h_0(x) = 0 + 0.5(2) = 1$$

$$h_0(x) = 0 + 0.5(3) = 1.5$$



minimum error at this point
or Global minima

we need to reach global minima so we cannot change the value of θ_0, θ_1 manually. there is some mechanism called convergence algorithm that optimize the change of θ value.

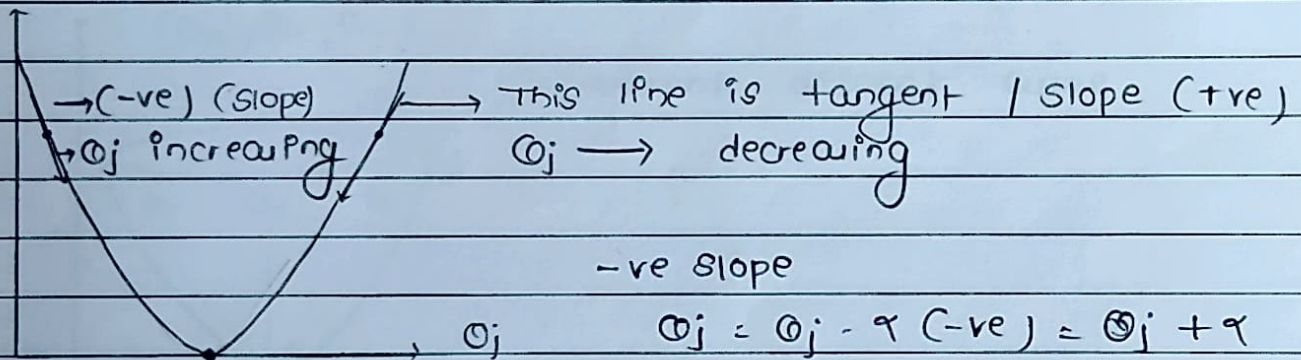
CONVERGENCE ALGORITHM :

Repeat until the convergence

$$\left\{ \begin{array}{l} \text{slope} \\ \uparrow \\ \theta_j = \theta_j - \alpha \frac{\partial J(\theta_j)}{\partial \theta_j} \end{array} \right\} \left\{ \begin{array}{l} \text{learning rate decides} \\ \text{the convergence speed} \end{array} \right\}$$

learning rate

$J(\theta_j)$



+ve slope

$$\theta_j = \theta_j - \alpha (+ve) = \theta_j - \alpha$$

α - learning rate

if α is very small \rightarrow will more time to reach global minima

if α is very large \rightarrow it will jump here & there won't reach global minima.

it should be around $\rightarrow 0.001$ for smaller steps.

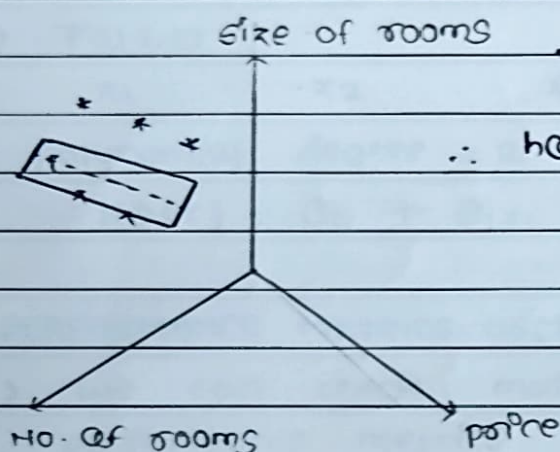
MULTIPLE LINEAR REGRESSION

> Two or more independent variable and one dependent variable.

> Example:

x_1	x_2	y
No. of rooms	Size of rooms	price

>



Try to find best fit plane

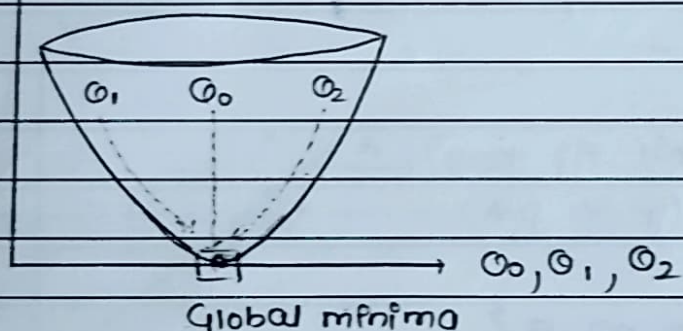
$$\therefore h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$\theta_1, \theta_2 \rightarrow$ slope or coefficient

$\theta_0 \rightarrow$ intercept

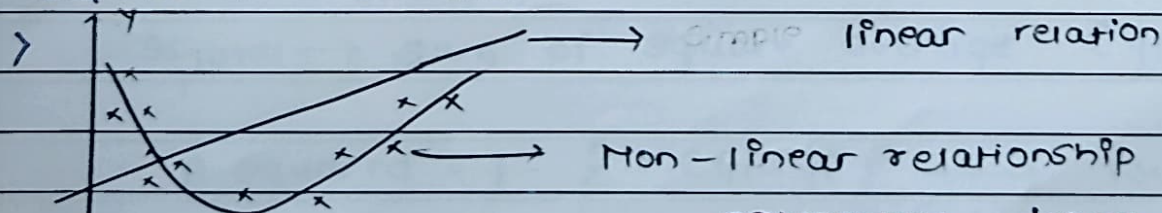
> $J(\theta_0, \theta_1, \theta_2)$

Gradient descent curve



POLYNOMIAL REGRESSION

> Non-linear relationship between independent and dependent variable



polynomial degree = n

$$h_{\theta}(x) = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n$$

FOR EDUCATIONAL USE

MULTIPLE POLYNOMIAL REGRESSION

> Multiple independent feature & one dependent feature

> Assumption

- linear relationship
- normally distributed
- no multicollinearity

> Dataset

x_1 x_2 x_3 $y \rightarrow$ o/p feature

polynomial degree = 2

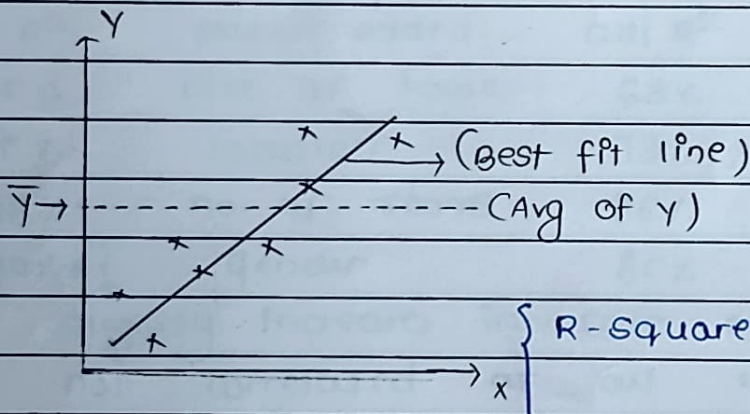
$$\therefore h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2$$

PERFORMANCE METRICS USED IN LINEAR REGRESSION

> We can check model is good or not using performance metrics.

1.) R-SQUARED

> Measure the performance of the model



$$\left\{ R\text{-Squared} = 1 - \frac{SS_{\text{Res}}}{SS_{\text{Total}}} \right\}$$

$SS_{\text{Res}} \rightarrow$ Sum of square residual $(y_i - \hat{y}_i)^2$

$SS_{\text{Total}} \rightarrow$ Sum of square average $(y_i - \bar{y})^2$

$$R\text{-Squared} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \rightarrow \text{low value}$$

$$\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \hat{y}_i)^2} \rightarrow \text{high value}$$

$$[R\text{-squared} \leq 1]$$

If $R\text{-squared} = 0.85 = 85\%$ accuracy
 $= 0.75 = 75\%$ accuracy

ADJUSTED R-SQUARED

Example: dataset

x_1	x_2	x_3	y
Size of house	No. of room	Location	price

> R^2 keep on increasing when we add more independent feature, whether the independent feature is highly correlated or not it will keep on increasing as the number of independent feature increasing.

R^2	feature added	adj R^2	independent feature
65%	Size of house	63%	$P=1$
75%	location	73%	$P=2$
88%	no. of rooms	86%	$P=3$
90%	Gender	85%	$P=4$

slightly increase in case of Gender which is not correlated at all with price.

$$\left\{ \text{Adjusted } R^2 = 1 - \frac{(1-R^2)(N-1)}{N-P-1} \right\}$$

$N \rightarrow$ Number of datapoint

$P \rightarrow$ Number of independent feature

Adjusted R^2 is best metrics to evaluate the model

TYPES OF COST FUNCTION

4.) MSE (Mean square error)

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

$y_i \rightarrow$ Actual, $\hat{y}_i \rightarrow$ predicted

$$\therefore \hat{y}_1 = \theta_0 + \theta_1 x$$

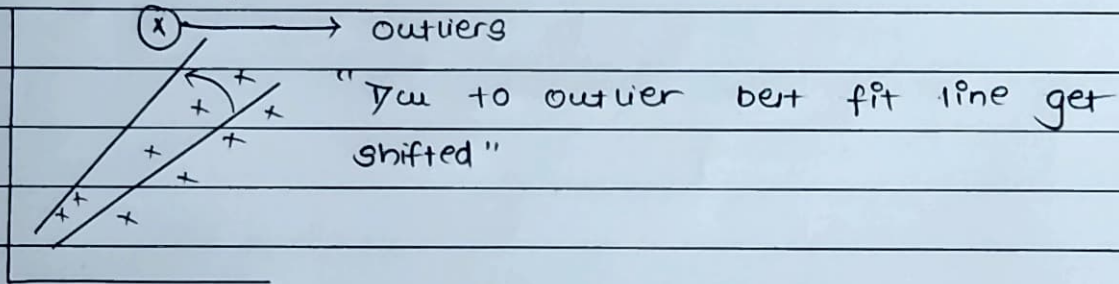
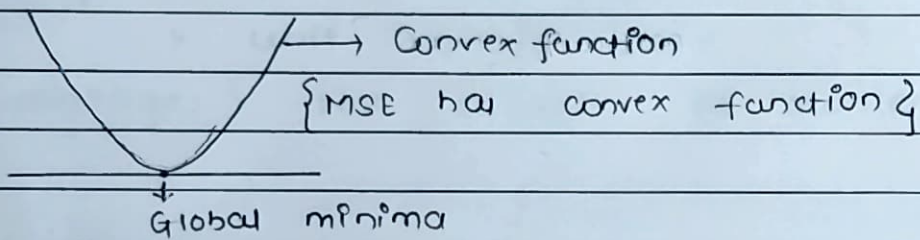
Advantage: } this equation is differentiable.

It also has one global minima.

Disadvantage: > Not robust to outliers

> not same unit

eg: $RS(INR) \longrightarrow (INR)^2$



Eg: Experience \rightarrow independent feature

Salary \rightarrow dependent feature

$(y - \hat{y})^2$ (laks)² \rightarrow unit changing \rightarrow time complexity
increase $\uparrow \uparrow$

we don't do scaling for dependent feature

TYPES OF COST FUNCTION

1.) MSE (Mean square error)

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

$y_i \rightarrow$ Actual, $\hat{y}_i \rightarrow$ predicted.

$$\therefore \hat{y}_i = \omega_0 + \omega_1 x$$

Advantage: > this equation is differentiable.

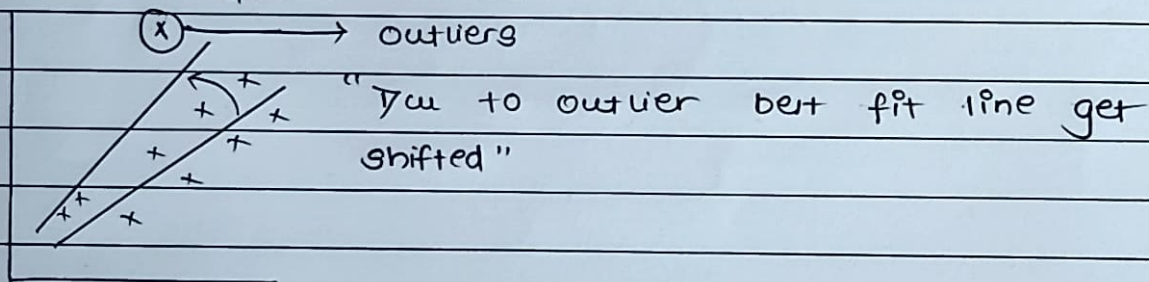
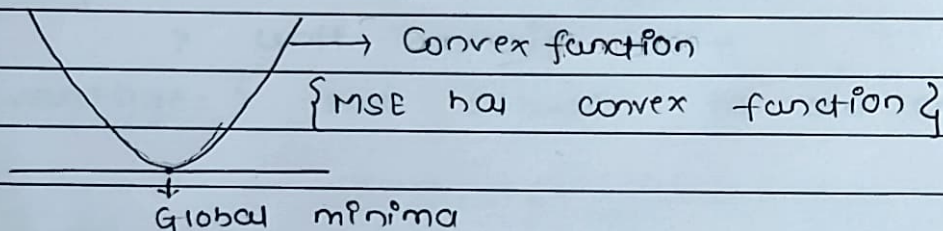
> it also has one global minima.

Disadvantage: > Not robust to outliers

> not same unit

eg: RS (INR) \rightarrow (INR)²

MSE has convex function



Eg: Experience \rightarrow independent feature

Salary \rightarrow dependent feature

$(y - \hat{y})^2$ (laks)² \rightarrow unit changing \rightarrow time complexity increases $\uparrow \uparrow$

ie don't do scaling for dependent feature

2.) MAE (Mean Absolute Error)

$$MAE = \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

Advantage: > robust to outliers
> same unit

Disadvantage: > convergence takes time
> optimization is complex

3.) RMSE (Root Mean Squared Error)

$$RMSE = \sqrt{MSE}$$

Advantage: > it is differentiable
> unit remain same

Disadvantage: > Not robust to outliers.