

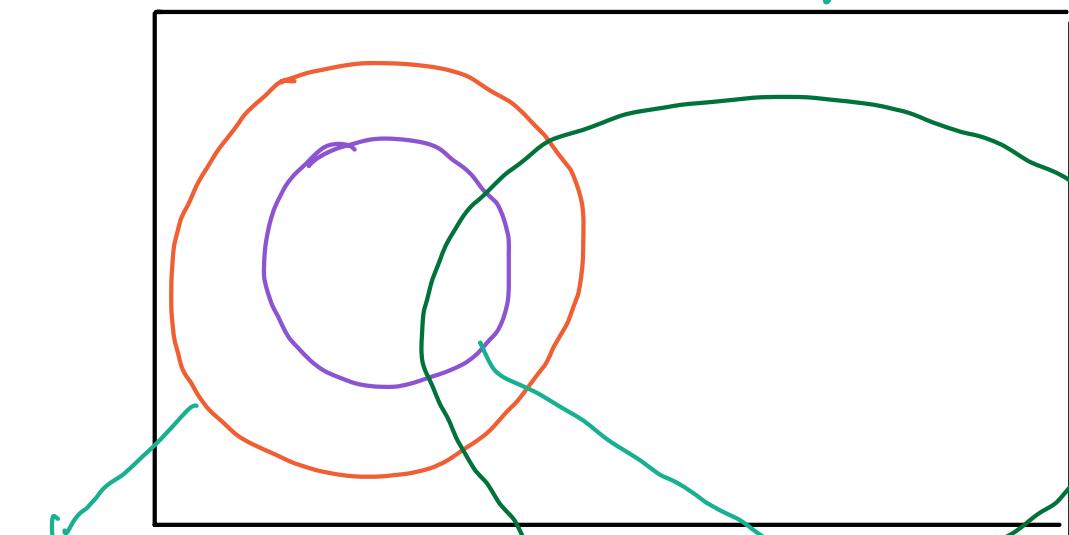
Universe of Data Science

- ① Artificial Intelligence
- ② Machine Learning
- ③ Deep Learning
- ④ Data science

AI → To create an application which can perform its own tasks without any human intervention

Eg: Netflix Recommendation system

Data science
↓
overlaps everything



ML

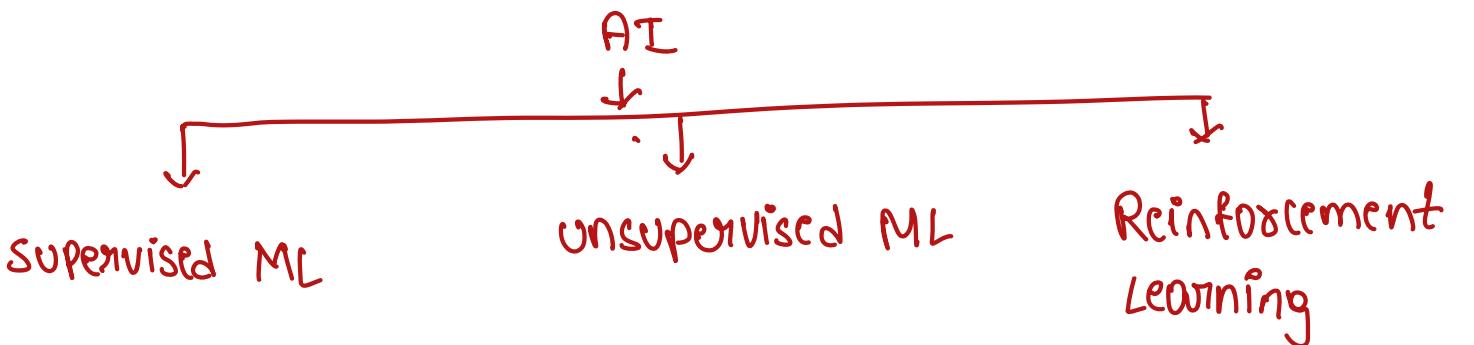
↓
It provides stats tools to analyse, visualize, predict, forecast the data

Deep Learning

↓
It was developed to mimic how humans learn

↓
multi layer neural network

3 Types of Machine Learning:



Supervised ML Technique

Regression

classification

→ Task is Predict House Price

Dataset

Independent features

Size of House

of Rooms

5000

5

6000

6

-

-

-

-

(i) Dependent or O/P feature

(ii) if output is continuous

↓
Regression Problem Statement

(iii) if output is discrete / categorical

↓
Classification Problem Statement

Classification Example

→ I/P (in) Dependent Feature

No. of Study Hours

No. of Play Hours

7

3

2

6

O/P (out) Dependent Feature

Pass / Fail

Pass

Fail

↓
categorical / Discrete

so ↓ classification

Dependent (or)
Output

Feature

Price

450 K

500 K

-

-

-

This value
Keeps on
changing
Based on

↓
Input (Independent Features)
Target feature is
continuous value

↓
Here

So it's a Regression
Problem Statement

Binary Classification (2 classes)

Multi Class Classification (> 2 classes)

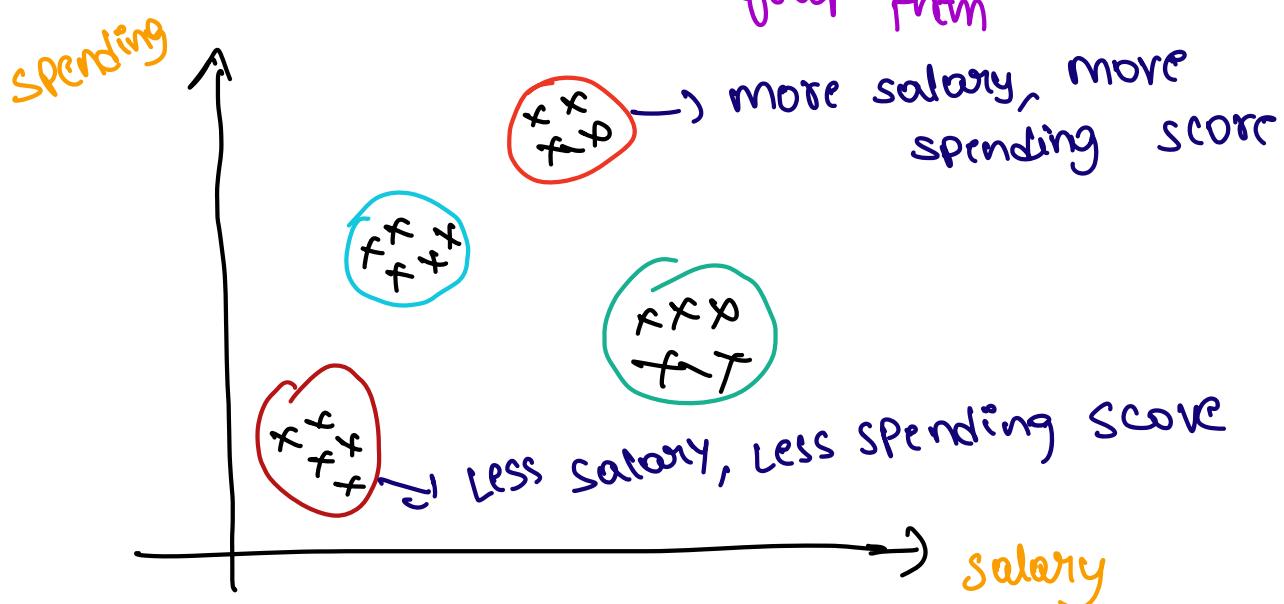
② Unsupervised ML Technique:

Eg: Customer Segmentation

- NO specific output / dependent Feature
- NO Labelling data

<u>salary</u>	<u>spending score (1-10)</u>	<u>E-commerce company</u>
20000	9	
45000	2	
-	-	↓
-	-	Email with some
-	-	discount coupons to
-	-	specific spending score

→ Here we create clusters and group them



Supervised ML Algorithms

- ① Linear Regression
 - ② Ridge & Lasso Regression
 - ③ Elastic Net
 - ④ Logistic Regression → classification Algorithm
 - ⑤ Decision Tree → Both Regression / classification
- } → Regression Algorithms

- ⑥ Random Forest
 - ⑦ Ada Boost
 - ⑧ Xg Boost
- } Both Regression / classification

Unsupervised ML:

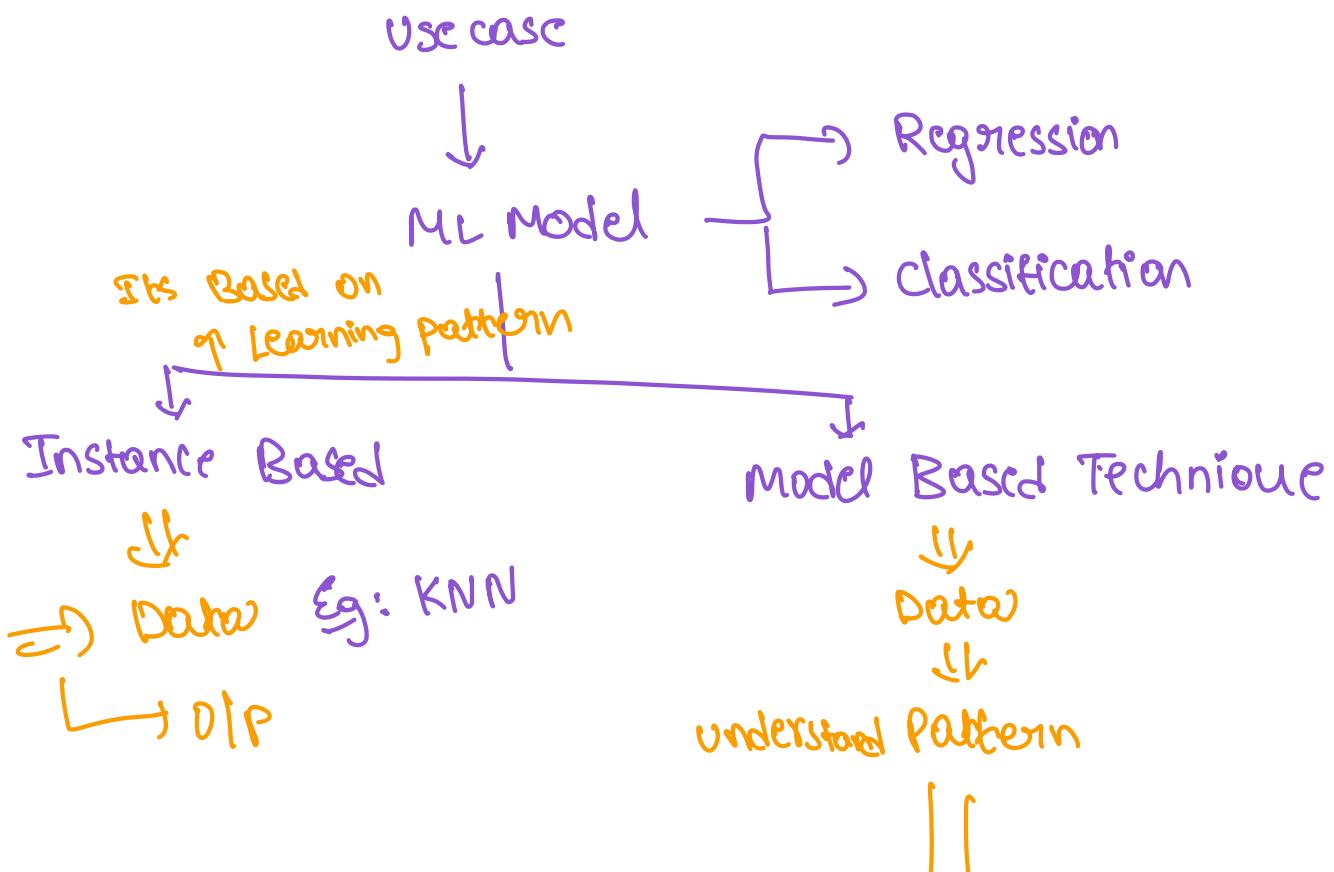
- ① K means
 - ② Hierarchical Mean
 - ③ DB scan
- } clustering Algorithms

Reinforcement Learning:

→ It will learn itself by Rewarding | Penality

Eg: Baby (0s) Kids learning

Instance Based Learning Vs Model Based learning





It will create Generalization method



whose Model can predict correctly for new data point

Model Based Learning



Instance Based Learning

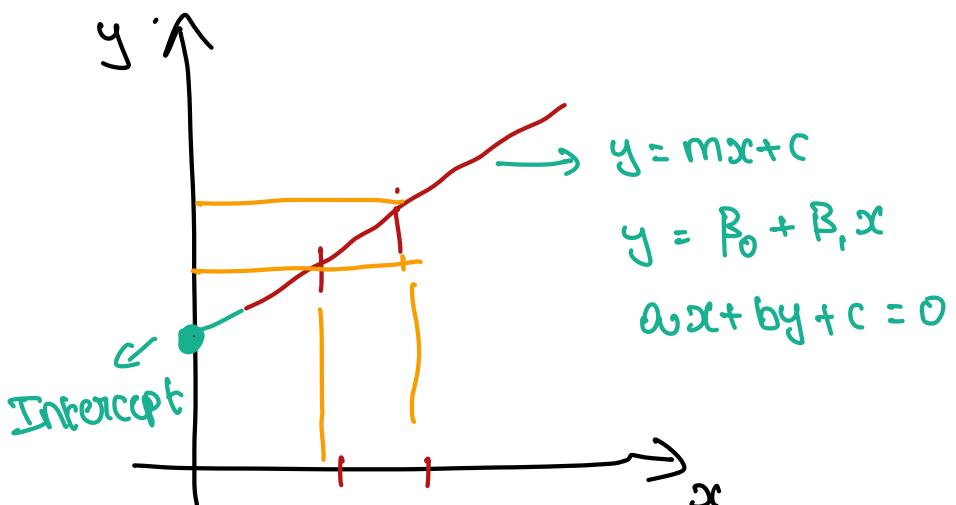


Usual/Conventional Machine Learning	Instance Based Learning
Prepare the data for model training ✓	Prepare the data for model training. No difference here ✓
Train model from training data to estimate model parameters i.e. discover patterns	Do not train model. Pattern discovery postponed until scoring query received
Store the model in suitable form (pickle format)	There is no model to store ✓
Generalize the rules in form of model, even before scoring instance is seen <i>Decision Boundary</i>	No generalization before scoring. Only generalize for each scoring instance individually as and when seen
Predict for unseen scoring instance using model	Predict for unseen scoring instance using training data directly
Can throw away input/training data after model training	Input/training data must be kept since each query uses part or full set of training observations
Requires a known model form	May not have explicit model form
Storing models generally requires less storage	Storing training data generally requires more storage
Scoring for new instance is generally fast	Scoring for new instance may be slow

Learning → memorizing → Instance Based
→ Generalizing → Model Based

Linear Algebra:

Equation of Line, 3d Plane and Hyperplane (n-Dimension)



① $y = mx + c \rightarrow$ Above line can be represented by this equation

$x, y \rightarrow$ co-ordinates

$m =$ Slope

In the unit movement in the x -axis, what is the movement in the y -axis

$c =$ Intercept

if $x = 0$ in the Equation $y = mx + c \Rightarrow y = mx + c$ $y = c$

where the straight line meets y -axis is called Intercept

② $ax + by + c = 0$

$$by + c = -ax$$

$$by = -ax - c$$

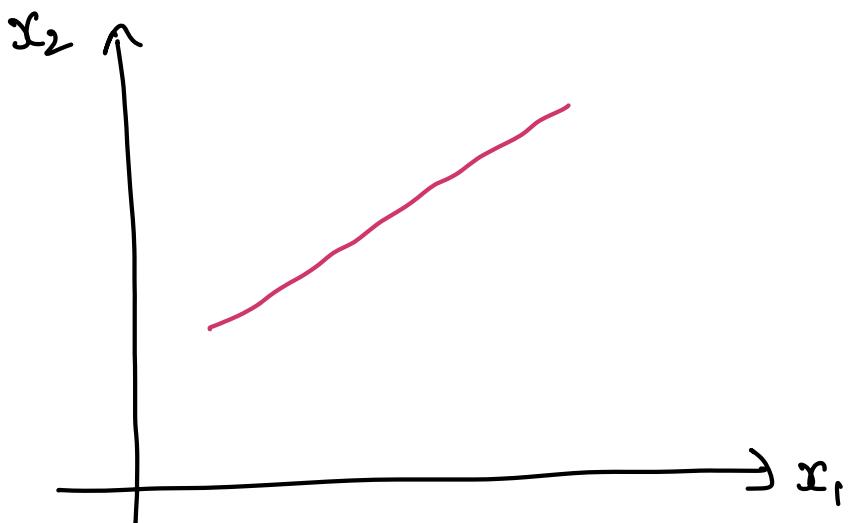
$$y = \left(\frac{-a}{b} \right) x - \frac{c}{b}$$

\downarrow

in the form of

$$\boxed{y = mx + c}$$

For Eg we have more than one dimension:



Above line can be written as:

$$\omega_1 x_1 + \omega_2 x_2 + b = 0$$

$$\omega_1 = \omega$$

$$\omega_2 = b$$

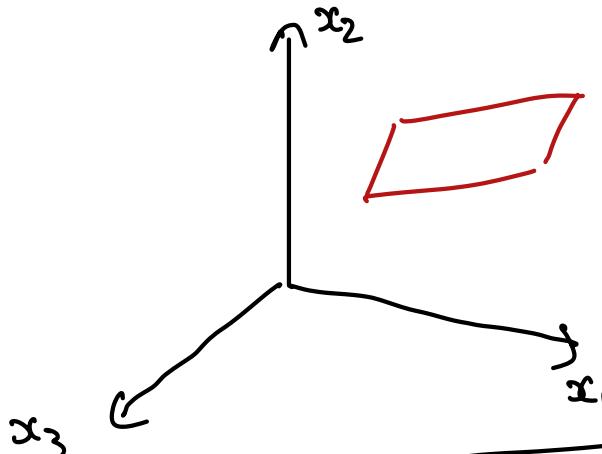
$$b = c$$

This can also be written as:

$$\omega^T x + b = 0$$

This is nothing but \Downarrow an equation of a straight line

For Eg what if we have 3 dimensions:



Here we don't draw
straight line
 \Downarrow
we draw plane

$$\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + b = 0$$

→ Above Equation can be written as:

$$w^T x + b = 0$$

$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$ $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

If we have to do the dot product of w and x we have to transpose one matrix

→ When we do matrix multiplication with ' w ' and ' x ' one of the matrix we have to transpose it
Why we have to transpose it?

(R1) In matrix multiplication, number of columns in the first matrix must equal the number of rows in the second matrix

(R2) Dot product Interpretation :

- Matrix multiplication involves taking the dot product of rows and columns of the matrices
- Transposing one matrix ensures that the dot product is taken between corresponding elements of the matrices, facilitating correct multiplication

What About n-dimensional Plane :

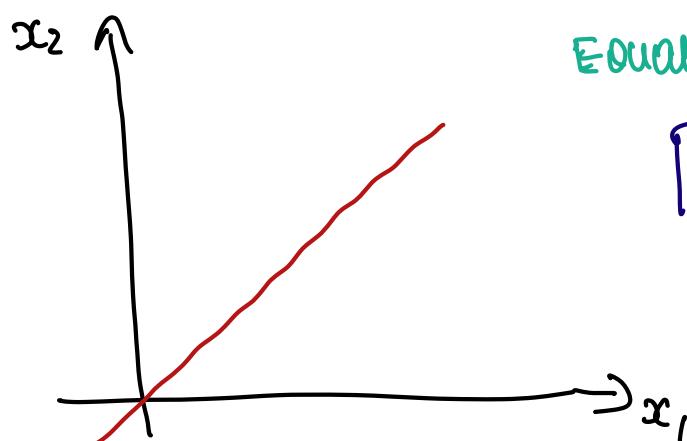
$$w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + b = 0$$



$$w^T x + b = 0$$

For Eg:

What if the straight line passes through origin?



Equation of a straight line:

$$w_1x_1 + w_2x_2 + b = 0$$



If straight line passes through origin
↓
means Intercept, $b = 0$

Equation of a straight line
passing through origin

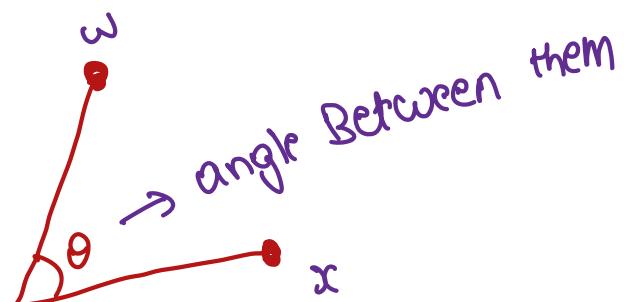
$$w_1x_1 + w_2x_2 = 0$$



$$w^T x = 0$$

Equation of a plane = $Tl_n = w^T x = 0$

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{bmatrix} \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$



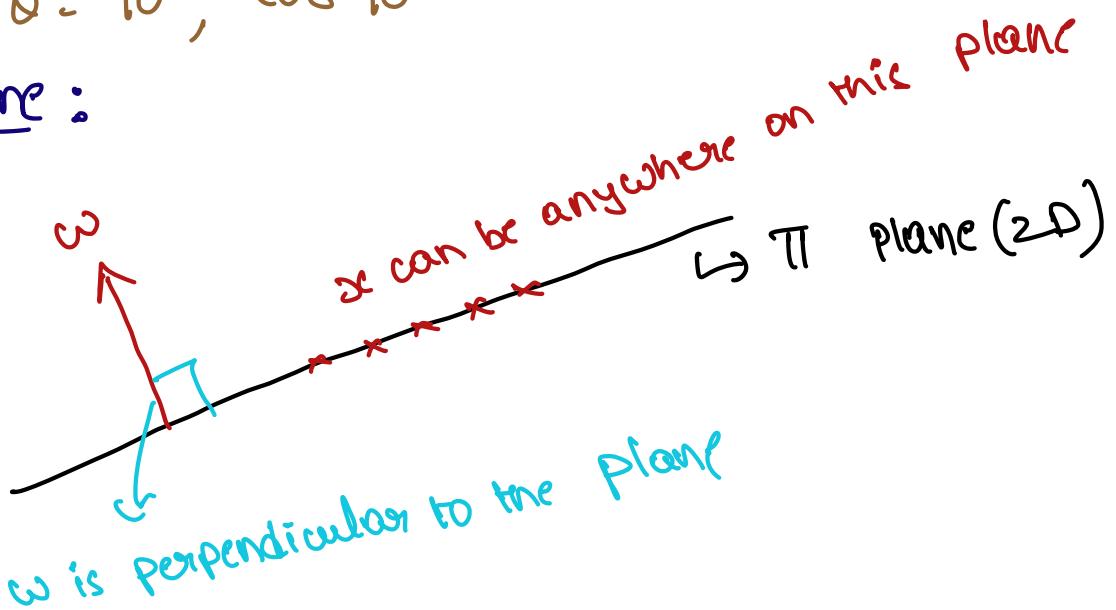
$$w \cdot x = w^T x = \|w\| \|x\| \cos \theta = 0$$

↓
magnitude of w * magnitude of x * Angle Between them

When $\cos \theta = 0$?

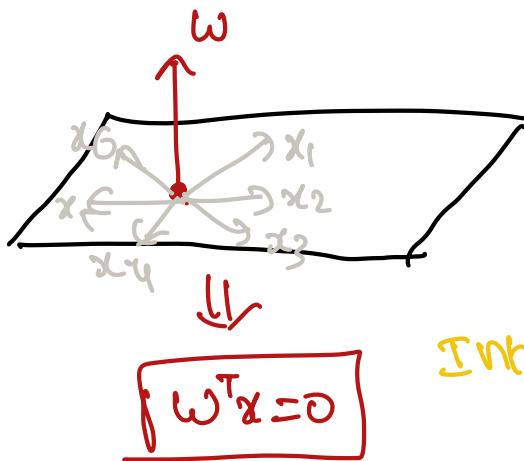
If $\theta = 90^\circ$, $\cos 90^\circ = 0$

In 2D plane:



In 3D plane:

' w ' is always perpendicular to the plane



Machine Learning Algorithms :

① Simple Linear Regression :

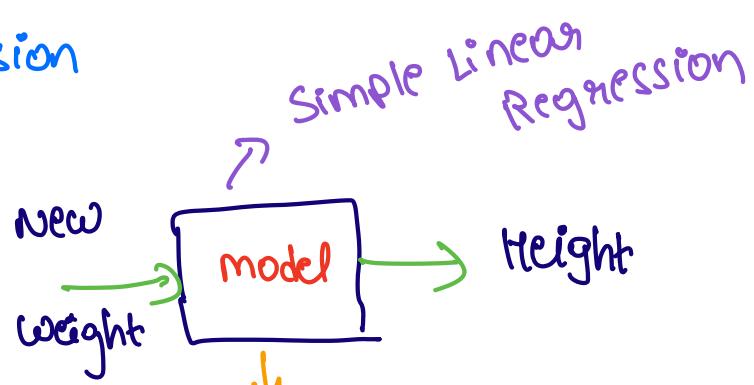
Supervised ML \rightarrow Regression

For eg :

Dataset

	weight	Height
Input /	74	170 cm
Independent Feature	80	180 cm
	75	175.5 cm
Feature	:	:

Output / Dependent Feature



We should train this model with Dataset

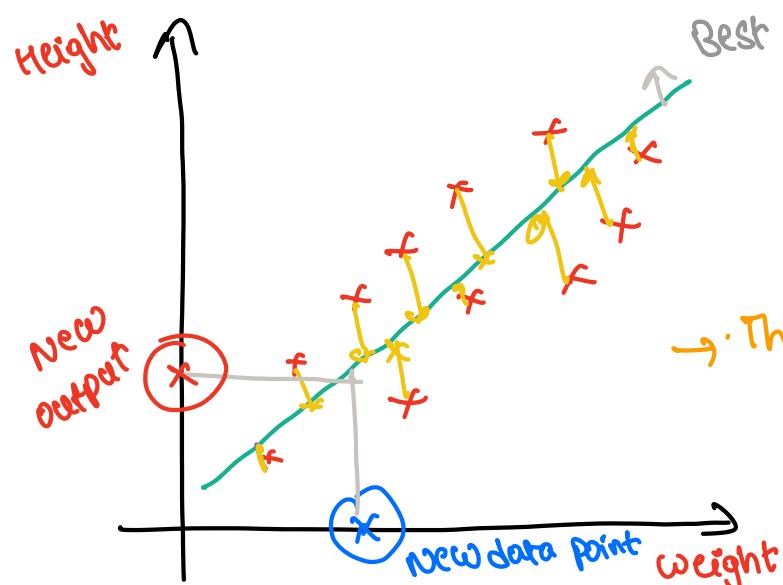
So that model can predict Height for New weight

which is not in Training set

Simple Linear Regression \rightarrow one Input feature

Multiple Linear Regression \rightarrow more than one Input Feature

What are we trying to do in Geometric way :



$x \rightarrow$ Input and output

data points (or)

observations

\rightarrow The summation of all distances
Should be minimum

That will be the Best Fit line

What are we trying to do in Linear Regression ?



Trying to fit a Best line among all the data points

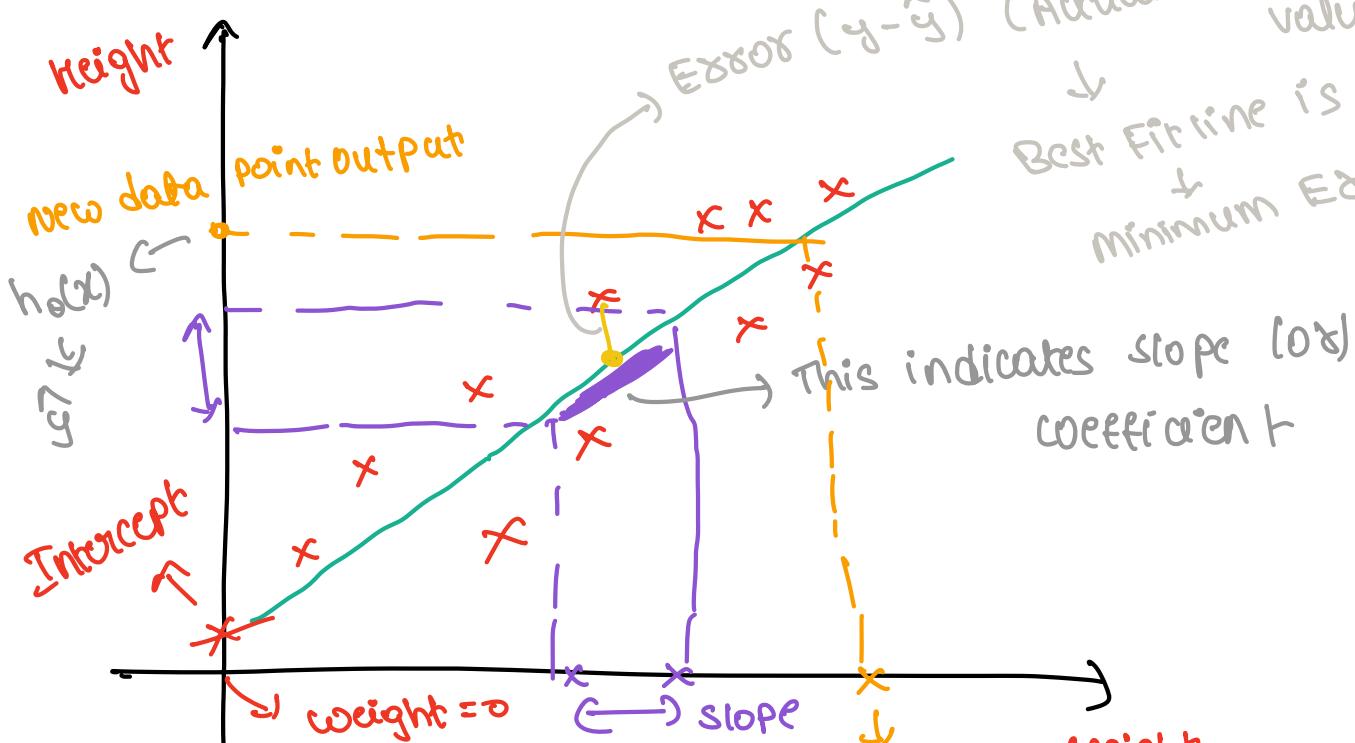


This Best Fit line will help us predictions for the new weight



All the distance between Actual value and Predicted value is called Error

Mathematical Equation of Best Fit Line :



Equation of Straight line :

$$y = mx + c$$

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

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

All are same

But mostly we will use this Equation

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$x \rightarrow$ independent (or) input Feature

$\theta_1 \downarrow$
It is weight in our case

$\theta_0 \rightarrow$ Intercept

$$\text{if } x=0 \rightarrow h_{\theta}(x) = \theta_0$$

\downarrow
when $x=0$, i.e., weight = 0

\downarrow
where our best fit line meet y-axis

$\theta_1 =$ Slope (or) Coefficient

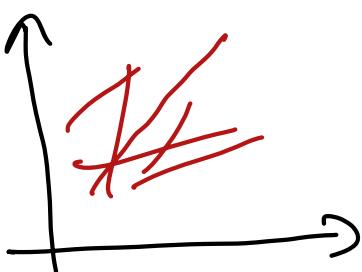
\downarrow
For a unit movement on x-axis, what is the movement w.r.t. y-axis

For n-Features :

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n$$

Predicted Output $\rightarrow h_{\theta}(x) \rightarrow \hat{y}$ (y hat)

How we are going to get Best Fit line ?



- ① Draw one simple straight
- ② Then we are going to change the angle and rotate it to find Best Fit line

$$h_0(x) = \theta_0 + \theta_1 x$$

θ_0 = Intercept

θ_1 = Slope (or) Coefficient

→ By changing Intercept (θ_0) and Slope (θ_1) we are going to rotate the straight fit line to find the Best Fit line.

Cost Function :

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

↓
Predicted points ↓
Actual (or) True points

$h_0(x) - y \rightarrow$ Gives Error Value



The reason why we are squaring ?



Because the cost function is "Mean Squared Error"

Final Aim , what we need to solve

↓↓↓

→ minimize cost function $J(\theta_0, \theta_1)$

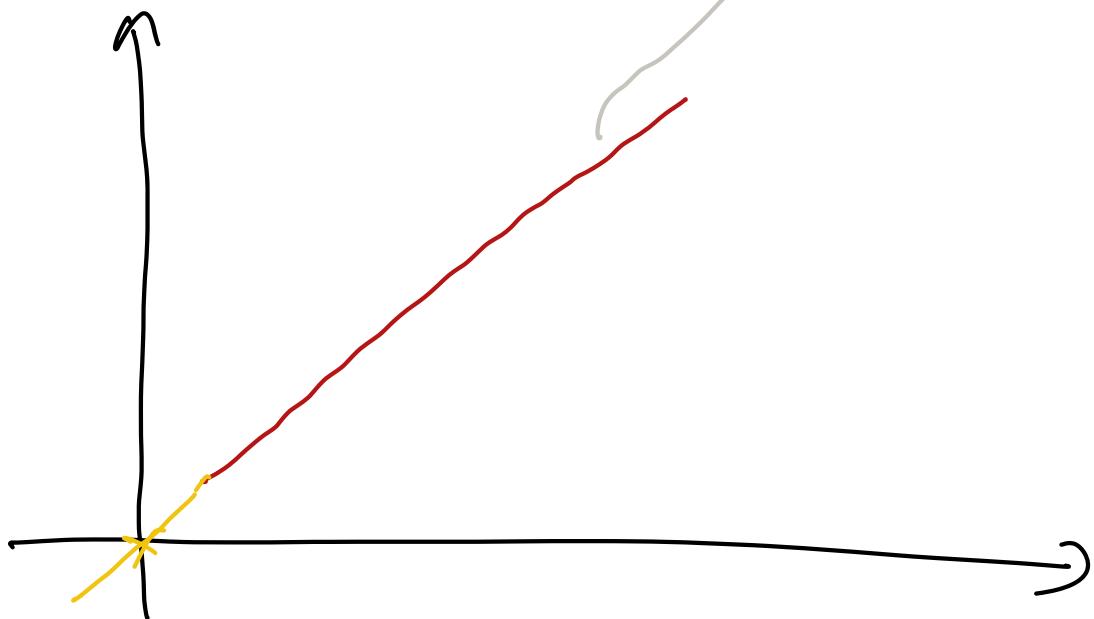
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

This value should be minimal by changing (θ_0, θ_1)

How we are going to do this?

Equation of straight line :

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



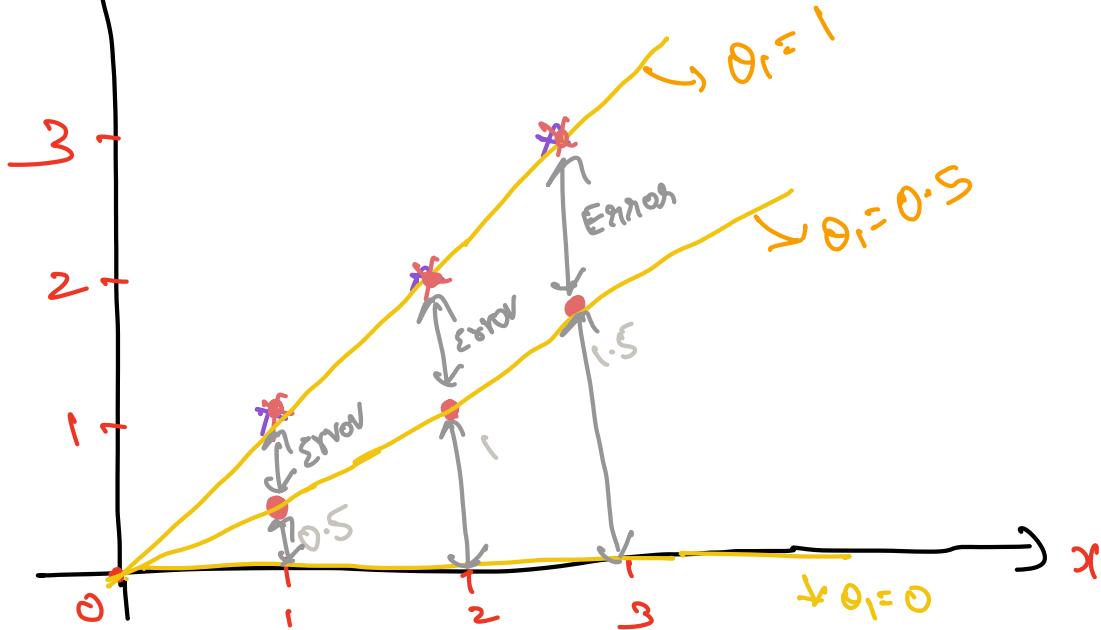
→ Assume our straight line passing through origin

$$\theta_0 = 0$$

$$h_{\theta}(x) = \theta_1 x$$

Dataset

x	y
1	1
2	2
3	3



\times → Actual point
 \bullet → Predicted point

line Equation : $h_{\theta}(x) = \theta_0 + \theta_1 x$

Let $\theta_1 = 1 \{ \text{slope} \}$ \downarrow we will keep on changing the slope

$$h_{\theta}(x) = 1, \text{ when } x = 1$$

$$h_{\theta}(x) = 2 \quad x = 2$$

$$h_{\theta}(x) = 3 \quad x = 3$$

$\rightarrow \theta_0 = 0$ (since line passing through origin, Intercept $\theta_0 = 0$)

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$$

$m = \text{no. of datapoints}$, in our case $m = 3$

$$= \frac{1}{2 \times 3} \left[(1-1)^2 + (2-2)^2 + (3-3)^2 \right] = 0$$

Cost Function is zero

If you see all predicted value & Actual values same in Above graph

Let $\theta_1 = 0.5$ ($h_{\theta}(x) = \theta_1 x$)

\downarrow

$$h_{\theta}(x) = 0.5 \text{ if } x = 1 \rightarrow 0.5 \times 1 = 0.5$$

$$x = 2 \rightarrow 0.5 \times 2 = 1$$

$$x = 3 \rightarrow 0.5 \times 3 = 1.5$$

$$h_{\theta}(x) = 1$$

$$h_{\theta}(x) = 1.5$$

Calculate Cost Function $J(\theta_0, \theta_1)$

$$\boxed{\theta_0 = 0}$$

Assumption, Line passing through origin

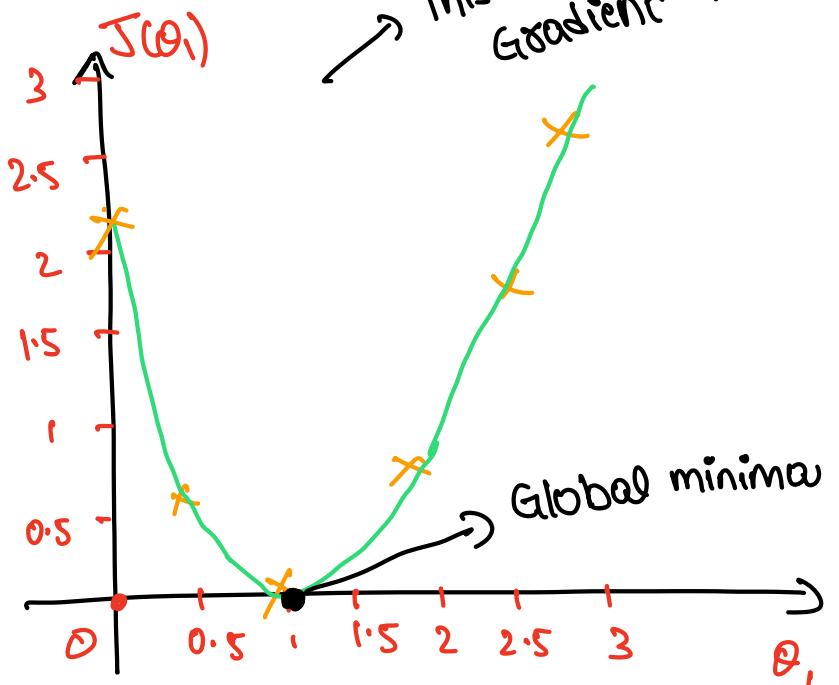
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m \left[h_\theta(x^{(i)}) - y^{(i)} \right]^2$$

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

$$\boxed{m=3}$$

$$\boxed{J(\theta_1) = \approx 0.58}$$

$$\text{when } \theta_1 = 0.5$$



$$\text{when } \theta_1 = 1, J(\theta_1) = 0$$

$$\text{when } \theta_1 = 0.5, J(\theta_1) = 0.58$$

$$\rightarrow \text{when } \theta_1 = 0$$

$$J(\theta_1) = \frac{1}{2 \times 3} \left[(0 - 1)^2 + (0 - 2)^2 + (0 - 3)^2 \right]$$

$$\approx 2.3 \rightarrow \text{Error quite high}$$

$$\text{when } \theta_1 = 0, J(\theta_1) = 2.3$$



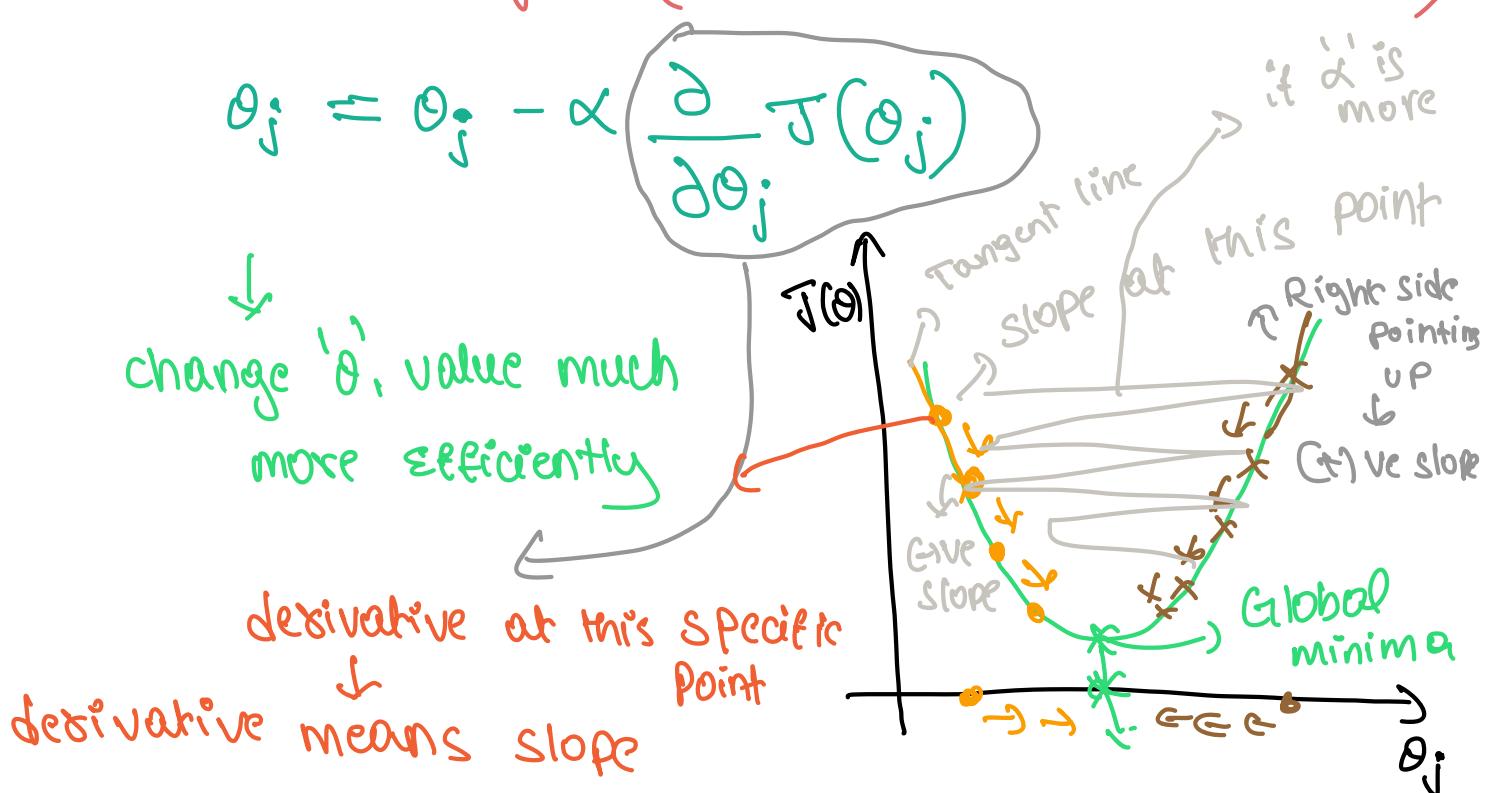
If we keep on changing
Above θ_1 values

↓
we will get Above
curve

- we have to ^{reach} Global minima point By changing (θ_0, θ_1)
i.e., we are reducing value of cost Function
- ↓
- We cannot keep on changing ' θ ' values practically
- ↓
- so, we will be needing convergence algorithm
- ↓
- We have to select one ' θ ' value and we have to find
a mechanism to select suitable ' θ ' values using
convergence algorithm

Convergence Algorithm ;

- Optimize the changes of ' θ ' value
- ↓
which is slope
- Repeat until convergence (means reach Global minimum)



How to find Derivative

→ Derivative is nothing but we are actually calculating slope



for calculating Derivative, we have to draw a tangent line at that point



we have to find whether it is positive slope (or) negative slope



so we will know whether we can Increase (or) decrease ' θ_j ' value

→ Just see the "right part" of Tangent line drawn at that point



If it is facing downwards



we will be calling it as negative slope

$$\rightarrow \theta_j = \theta_j - \alpha \frac{\partial J(\theta_j)}{\partial \theta_j}$$

From Above

↓
slope is (-ve)

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

(- \times - = +)

$$\theta_j = \theta_j + (+ve) \rightarrow \text{so we are increasing } \theta_j \text{ value}$$

→ How much θ_j value increase?



we don't know → This value will be decided by convergence algorithm by iterating

→ when slope is (+)ve

$$\theta_j = \theta_j - \alpha \left(\frac{\partial J(\theta)}{\partial \theta_j} \right) \rightarrow (+) \text{ ve slope}$$

$$\theta_j = \theta_j - \alpha (+\text{ve slope value})$$



θ_j will decrease



move towards Global Minima

α : Learning rate ($\alpha = 0.001$) → start with small value

→ Learning rate controls the speed at which convergence should happen

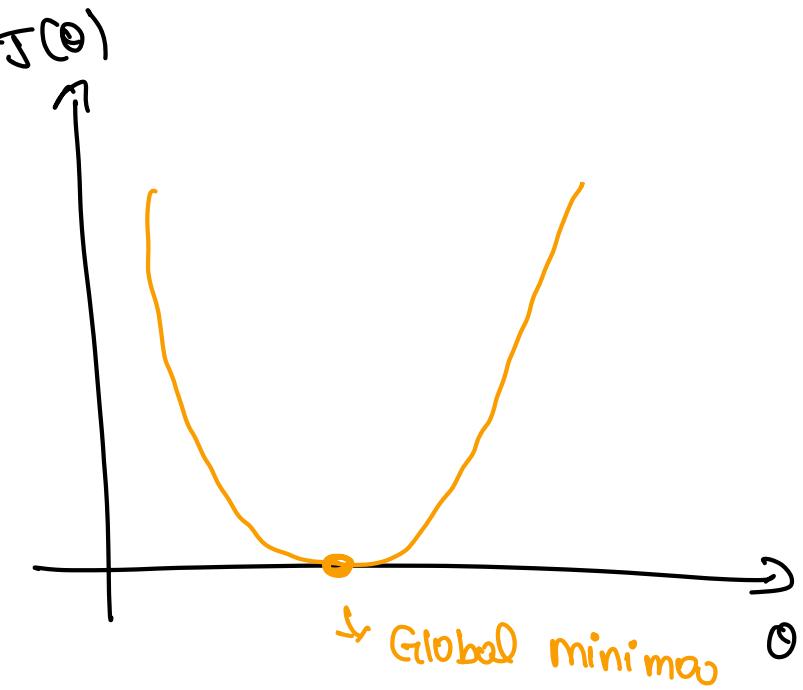
$\alpha \rightarrow$ Small \rightarrow model will take more time to converge

$\alpha \rightarrow$ Large \rightarrow model will converge Fastly

It may move very Fastly randomly and may not converge

Final Conclusion :

Gradient Descent :



Our Assumption :

→ Our Best Fit line passes through origin

i.e., $\boxed{\theta_0 = 0}$

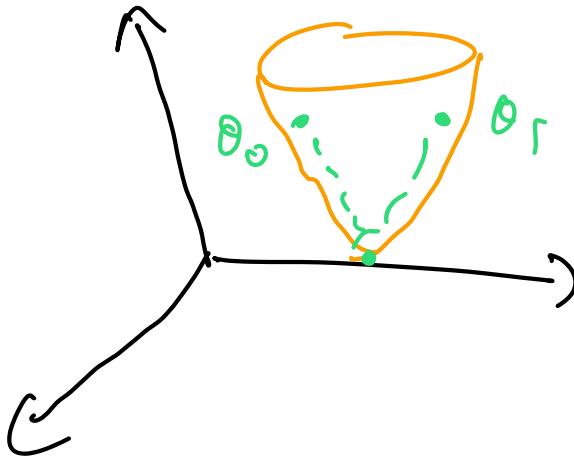
→ we only focussed on $\boxed{\theta_1 = \text{slope}}$ → changeable

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

This is done to visualize Gradient descent in
'2D' Diagram

For Eg: $\theta_0 \neq 0$ (i.e., our line not passing
through origin)

Gradient Descent in '3D'



convergence Algorithm :

repeat until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

Including θ_0

$j = 0 \text{ and } 1$

COST Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

→ if $j=0$



$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_0} \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \right]$$



$$2x^2$$

$$\frac{\partial}{\partial x} (x)^2 = 2x$$



$$\frac{\partial}{\partial x} (x)^n = n \cdot x^{n-1}$$

$$= \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(x)^{(i)} - y^{(i)} \right) x^{(i)}$$

$$h_{\theta}(x) = \theta_0 + \boxed{\theta_1 x} \rightarrow \text{constant}$$

→ if $j=1$

↓

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_1} \frac{1}{2m} \left[\sum_{i=1}^m \left[(\theta_0 + \theta_1 x^{(i)}) - y^{(i)} \right]^2 \right]$$

$$= \frac{1}{m} \sum_{i=1}^m \left((\theta_0 + \theta_1 x^{(i)}) - y^{(i)} \right) x^{(i)}$$

$\frac{\partial}{\partial \theta_1} [\theta_0 + \theta_1 x]$
↓
constant

derivative will be x

Repeat until convergence

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(x)^{(i)} - y^{(i)} \right) \\ \theta_1 := \theta_1 + \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(x)^{(i)} - y^{(i)} \right) x^{(i)} \end{array} \right.$$

multiple Linear Regression :

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

x = input Feature

θ_0 = Intercept

θ_1 = Slope

we won't change ' x '



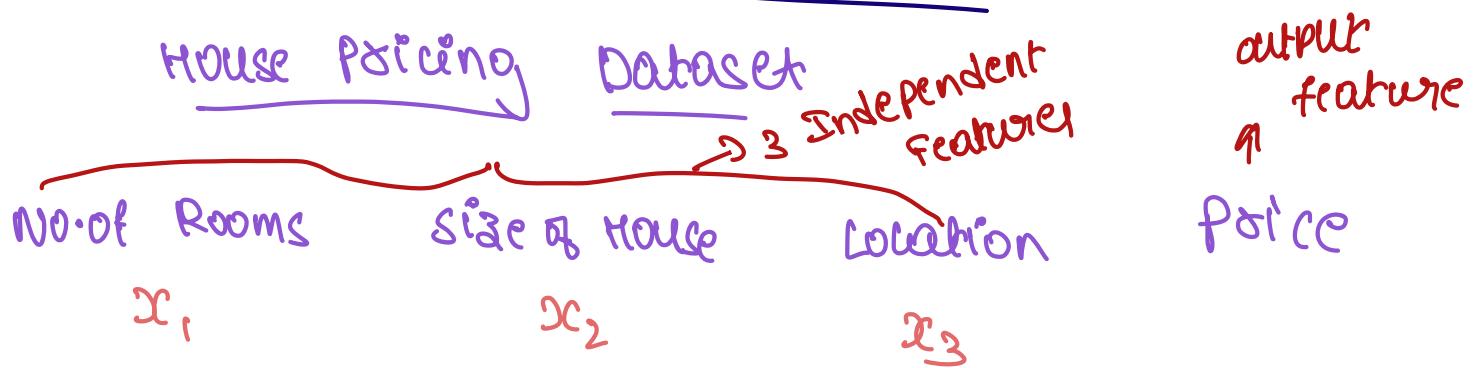
we only change

θ_0, θ_1



Find Best Fit Line

what if we have multiple Input Features:



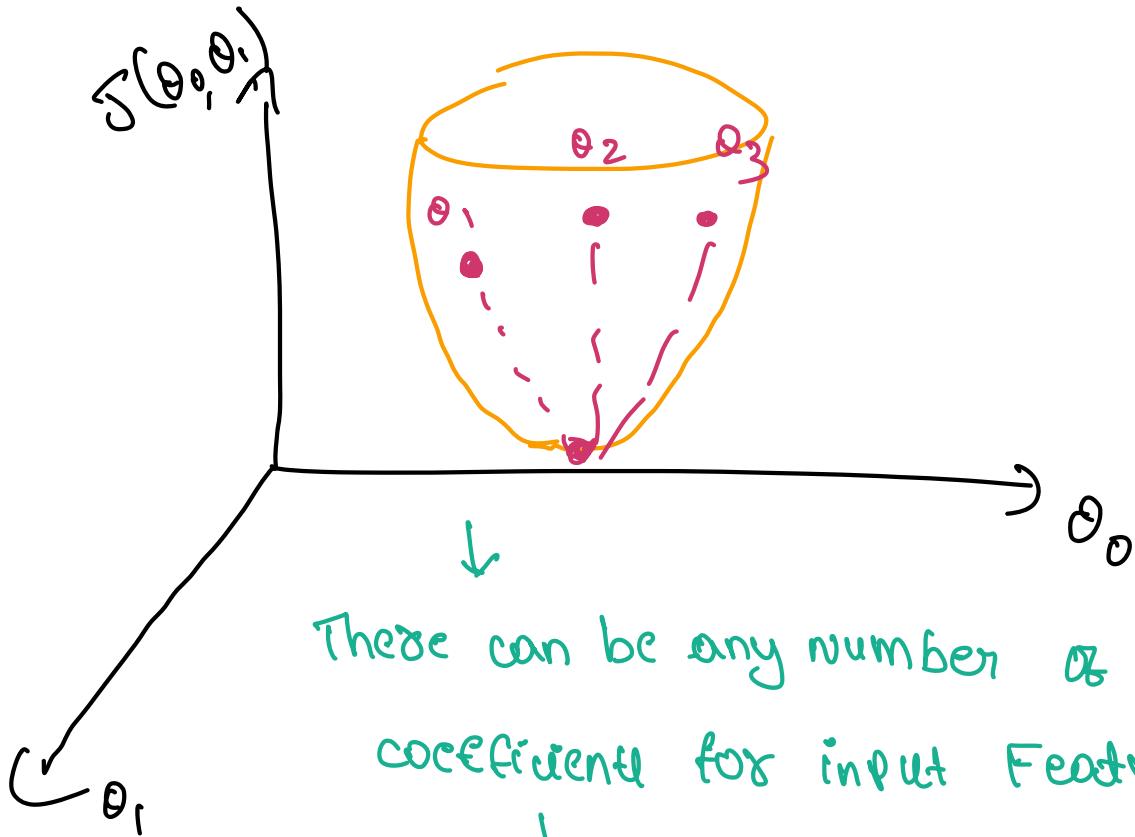
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$



Equation for multiple Linear Regression

$\theta_1, \theta_2, \theta_3$ = coefficient (or) slope of x_1, x_2, x_3

θ_0 = Intercept will be only one



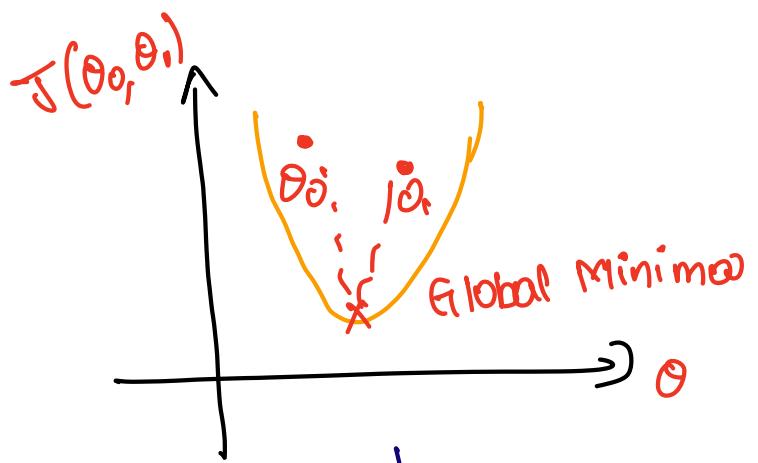
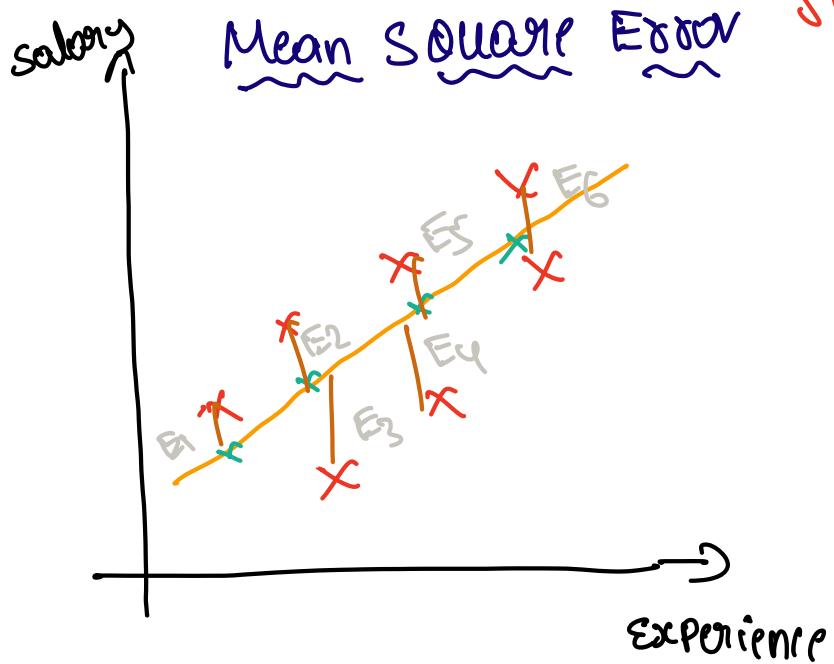
Their only aim is to reach Global minima

[Cost Function] \rightarrow Performance Metrics:

MSE \rightarrow Mean Square Error

MAE \rightarrow Mean Absolute Error

RMSE \rightarrow



This is for Gradient Descent

Data set

Input	Experience	Salary	O/P
-	-	-	
-	-	-	
-	-	-	
-	-	-	

→ we need to reduce the Cost Function

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

→ cost function

cost function, MSE = $E_1 + E_2 + E_3 + E_4 + E_5 + E_6$

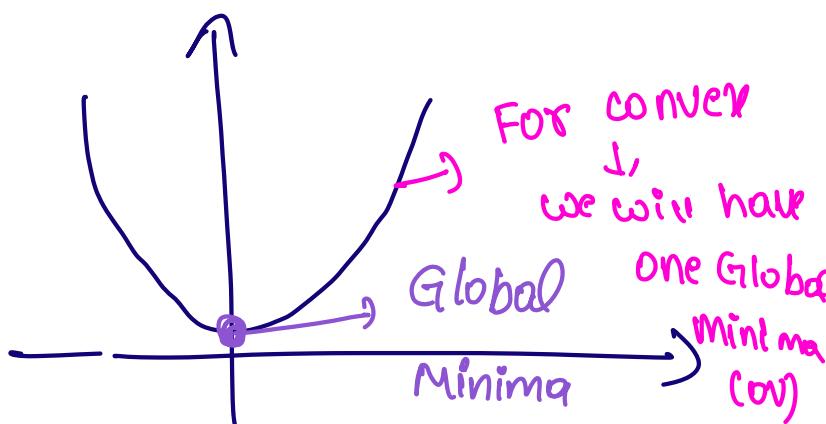
This value ↓ should be minimum

→ $(y - \hat{y})^2$ → It is Basically a Quadratic Equation

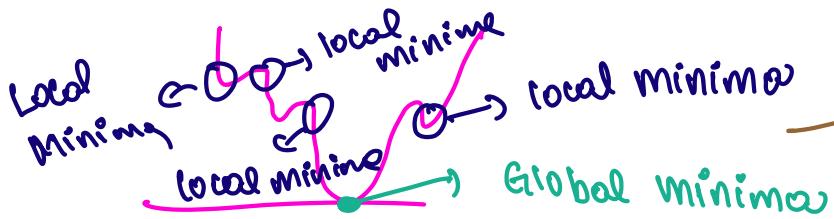
$$(a^2 - b)^2 = a^2 - 2ab + b^2$$

↓
Graph for Quadratic Equation

This is very
similar to
Gradient Descent



→ So By using (MSE) cost function we can easily reach Global Minima

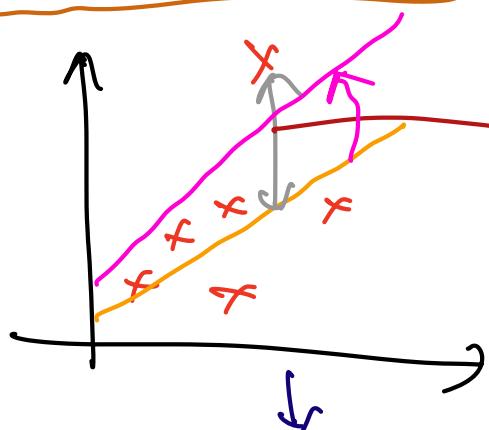


For non-convex Graph like this we will have many local minima

Advantages of MSE

- ① It is Differentiable at every point
- ② It has one local minima & global minima
- ③ Converge Faster

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



Disadvantages of MSE

- ① Not robust to Outlier.
- ② It is no longer in the same unit

Here we are penalizing the outlier
 ↓
 Error will Increase

MSE will penalize this outlier
 ↓
 our Best Fit line will move towards outlier (↗)

This is the Impact of outlier

↓
 MSE is Increasing

- At local minima, the convergence will get stuck, Because slope at local minima will be "Zero"
- So if we use any equations other than quadratic equations may have many local minima
- So it is Advisable to use quadratic equation

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

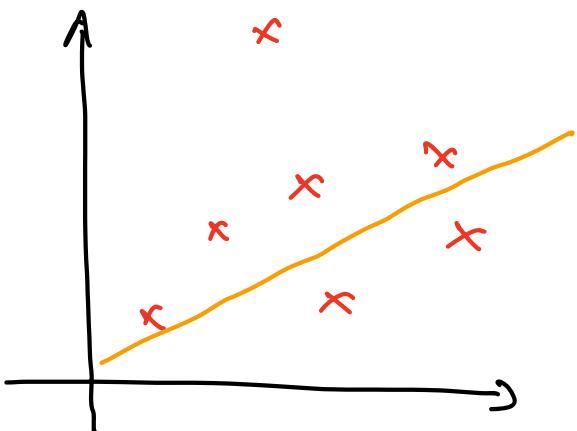
If unit of \hat{y} is lakhs

Doing $(y - \hat{y})^2$ → Becomes $(\text{lakhs})^2$

which is changing the unit

Mean Absolute Error (MAE)

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



- Here we are NOT squaring
- so it's robust to outliers
- It will be in the same unit

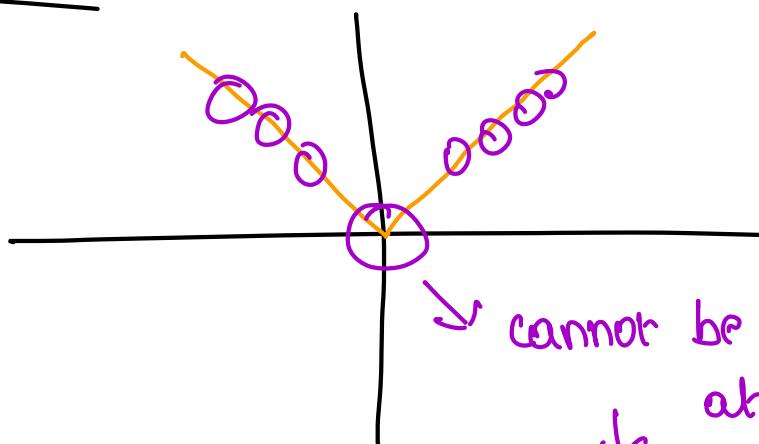
Advantages

- ① Robust to outliers
- ② It will be in the same unit

Disadvantages

- ① Convergence usually take more time
- ② Optimization is a complex task
- ③ Time consuming

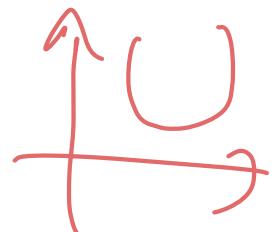
Graph of MAE :



↓ cannot be differentiable
at 'zero'
↓
So we create sub gradient
↓
Positive at Global Minima

RMSE [Root mean Square Error]

$$\begin{aligned} \text{RMSE} &= \sqrt{\text{MSE}} \\ &= \sqrt{\frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2} \end{aligned}$$



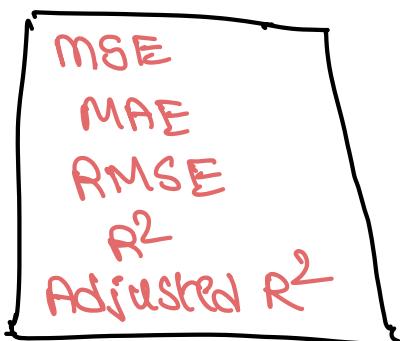
Advantage

- Same unit
- Differentiable

Disadvantage

- Not Robust to outliers

Regression Performance Metrics



R-Squared :

$$R_{\text{Squared}} = 1 - \frac{SS_{\text{Res}}} {SS_{\text{Total}}}$$

$$= 1 - \frac{\text{sum of squares of Residual}}{\text{sum of squares of Total}}$$

$$= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

$$\begin{aligned} &= 1 - \frac{\text{smaller number}}{\text{bigger number}} \\ &= 1 - \text{smaller number} \\ &\approx 1 \end{aligned}$$

For eg, $R_{\text{Squared}} = 0.70$, our model had 70%

Adjusted R² :

Dataset

size of house ↑ ↑ price ↑

Bed rooms, Location, Gender

accuracy

Both are having good positive correlation

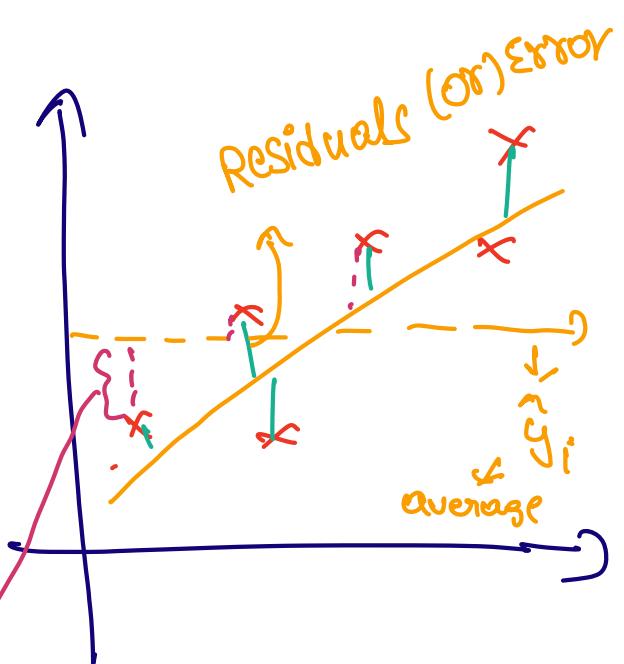
$R_{\text{Squared}} = 75\% \Rightarrow 0.75$

Bed rooms ∝ Price

$R_{\text{Squared}} = 80\% \Rightarrow 0.80$

Location ∝ Price

$R_{\text{Squared}} = 85\% \Rightarrow 0.85$



→ Gender doesn't have any correlation with price

$$R^2 \text{ squared} = 87\% \Rightarrow 0.87$$

→ AS features Increases, Due to the formula of R square value Increases

Even though the input feature is NOT at all correlated with output feature

↓
R square value Increases

→ Adjusted R square will resolve this problem

↓

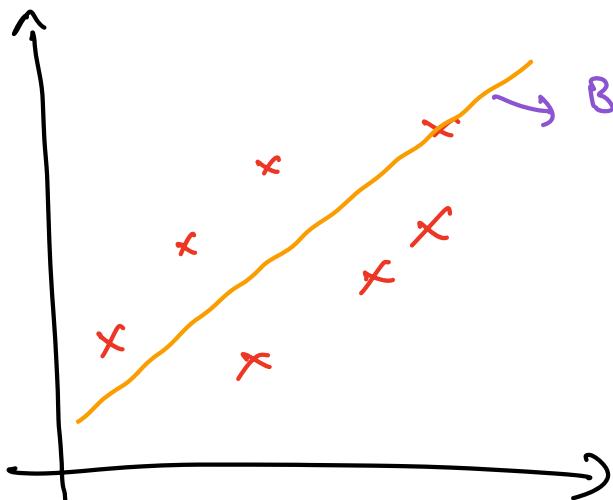
By penalizing the features which are not correlated with o/p features

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

N = no. of data points

p = no. of Independent features

Linear Regression using OLS (Ordinary Least Square)



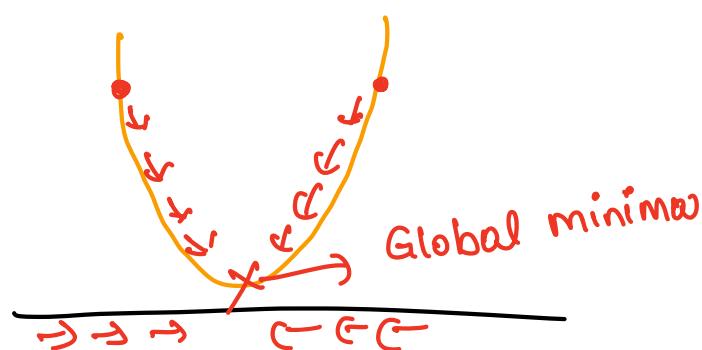
Best line was found
using Optimized Algorithm Before

↓
Gradient Descent

↓
Loss Function

↓
Every point we calculated
Derivative

↓
Drive towards Global
minimum



OLS Technique :

- using formula and calculate β_0 & β_1
- Aim of OLS is to reduce cost function error

ordinary Least Square :

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

↓
Aim is to calculate β_0 & β_1

calculative derivatives w.r.t. β_0

$$\frac{\partial S}{\partial \beta_0} (\beta_0, \beta_1) = \frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (0 - 1 - 0)$$

$$= -\frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0 \rightarrow ①$$

$$\frac{\partial \mathcal{L}^2}{\partial x} = 2x$$

$$\frac{\partial S}{\partial \beta_1} (\beta_0, \beta_1) = \frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (0 - 0 - x_i)$$

$$= -\frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (x_i) = 0 \rightarrow \textcircled{2}$$

EQ $\rightarrow \textcircled{1}$

$$\frac{-2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$-\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$-\sum_{i=1}^n y_i + n * \beta_0 + \beta_1 \sum_{i=1}^n x_i = 0$$

$$n \beta_0 + \beta_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$

$$\boxed{\beta_0 = \frac{\sum_{i=1}^n y_i}{n} - \beta_1 \frac{\sum_{i=1}^n x_i}{n}}$$

Intercept
can be
calculated by
using this
formula

$$\boxed{\beta_0 = \bar{y} - \beta_1 \bar{x}}$$

EQ $\rightarrow \textcircled{2}$

$$-\frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (x_i) = 0$$

$$-\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (x_i) = 0$$

$$\sum_{i=1}^n x_i y_i - \beta_0 \sum_{i=1}^n (x_i) - \beta_1 \sum_{i=1}^n (x_i)^2 = 0$$

$$\sum_{i=1}^n (x_i y_i - \beta_0 (x_i) + \beta_1 x_i^2) = 0$$

Replace $\beta_0 = \bar{y} - \beta_1 \bar{x}$

$$\sum_{i=1}^n (x_i y_i - (\bar{y} - \beta_1 \bar{x}) x_i - \beta_1 x_i^2) = 0$$

$$\sum_{i=1}^n (x_i y_i - \bar{x} \bar{y} + \beta_1 \bar{x} x_i - \beta_1 x_i^2) = 0$$

$$\sum_{i=1}^n (y_i - \bar{y} + \beta_1 \bar{x} - \beta_1 x_i) = 0$$

$$\sum_{i=1}^n \left[(y_i - \bar{y}) + \beta_1 (\bar{x} - x_i) \right] = 0$$

$$\sum_{i=1}^n (y_i - \bar{y}) + \sum_{i=1}^n \beta_1 (\bar{x} - x_i) = 0$$

$$\sum_{i=1}^n \beta_1 (\bar{x} - x_i) = - \sum_{i=1}^n (y_i - \bar{y})$$

$$\beta_1 = \frac{- \sum_{i=1}^n (y_i - \bar{y})}{\sum_{i=1}^n (\bar{x} - x_i)}$$

$$\beta_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})}$$

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

→ which is our slope
(or) coefficient

OLS Formulas

Intercept

$$\downarrow$$

$$\boxed{\beta_0 = \bar{y} - \beta_1 \bar{x}}$$

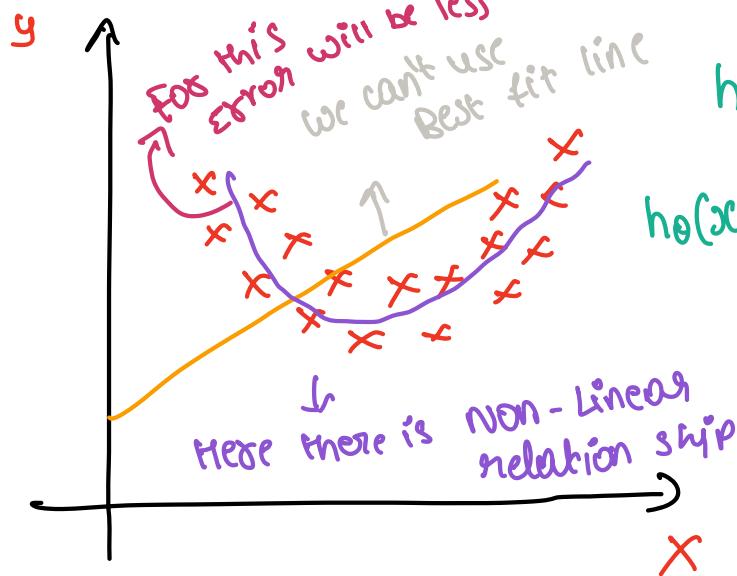
Coefficient

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

$$\begin{array}{c} x \\ \vdots \\ \bar{x} \end{array} \quad \begin{array}{c} y \\ \vdots \\ \bar{y} \end{array} \quad \begin{array}{c} (y_i - \bar{y}) \\ \vdots \\ \vdots \end{array} \quad \begin{array}{c} (x_i - \bar{x}) \\ \vdots \\ \vdots \end{array}$$

Polynomial Regression :

only on Independent Feature



$$h_0(x) = \beta_0 + \beta_1 x \rightarrow \text{Simple Linear Regression}$$

$$h_0(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

↓
Multiple Linear Regression

↓
more than one Independent Feature

what if my dataset is not Linear ? (Like Above)

→ In Both Simple Linear Regression, Multiple Linear Regression we assume Data points are Linear



In this case we will use polynomial Regression

polynomial Degrees

degree = 0



Simple Polynomial Regression (one IP, one OP)

if degree = 0



$$h_0(x) = \beta_0 \times x^0 \Rightarrow \text{constant value}$$

Polynomial degree = 1,

$$h_0(x) = \beta_0 \times x^0 + \beta_1 x^1$$



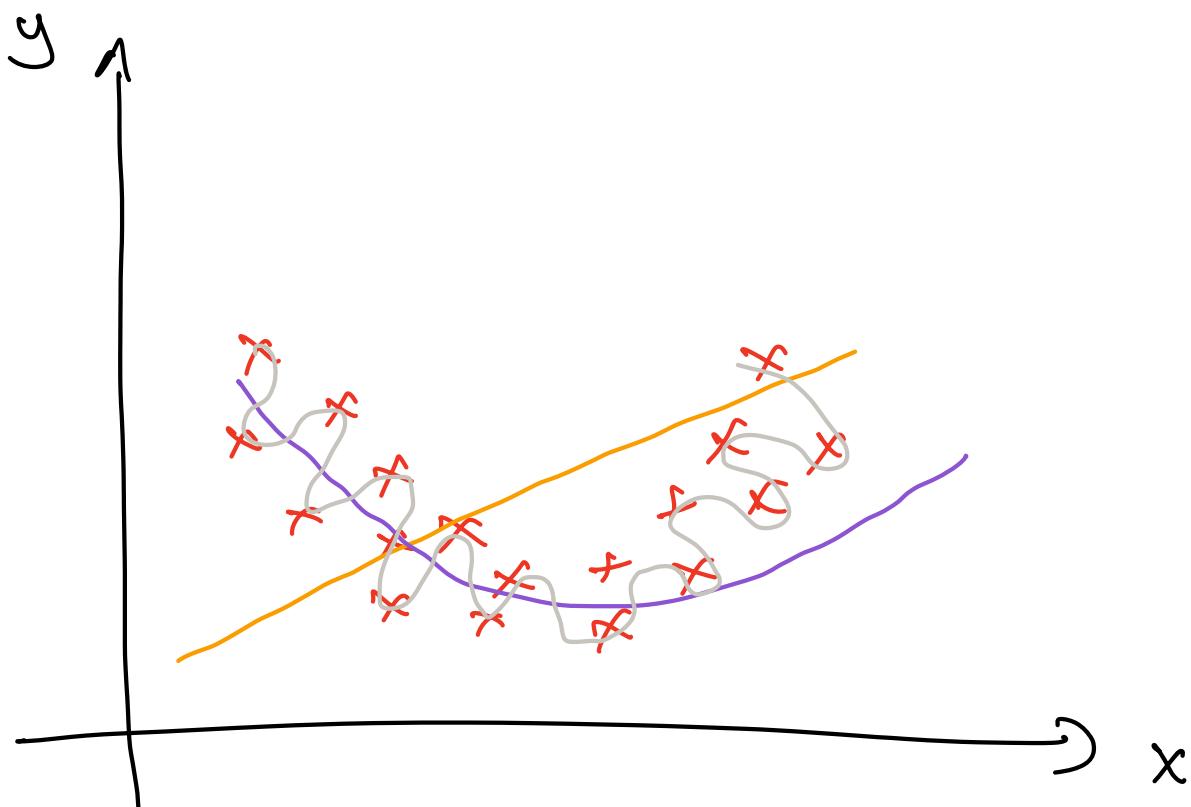
Simple Linear Regression

Polynomial degree = 2,

$$h_0(x) = \beta_0 x^0 + \beta_1 x^{(1)} + \beta_2 x^{(2)}$$

if degree = n,

$$h_0(x) = \beta_0 x^0 + \beta_1 x^{(1)} + \beta_2 x^{(2)} + \beta_3 x^{(3)} + \dots + \beta_n x^{(n)}$$



- As polynomial degree increases, our line (or) curve will fit more data points
- But if we are using much bigger Polynomial degree, the line will try to fit every other data point and becomes Overfit on data points
- So we have to find the Best degree for the curve to fit all data points without overfitting

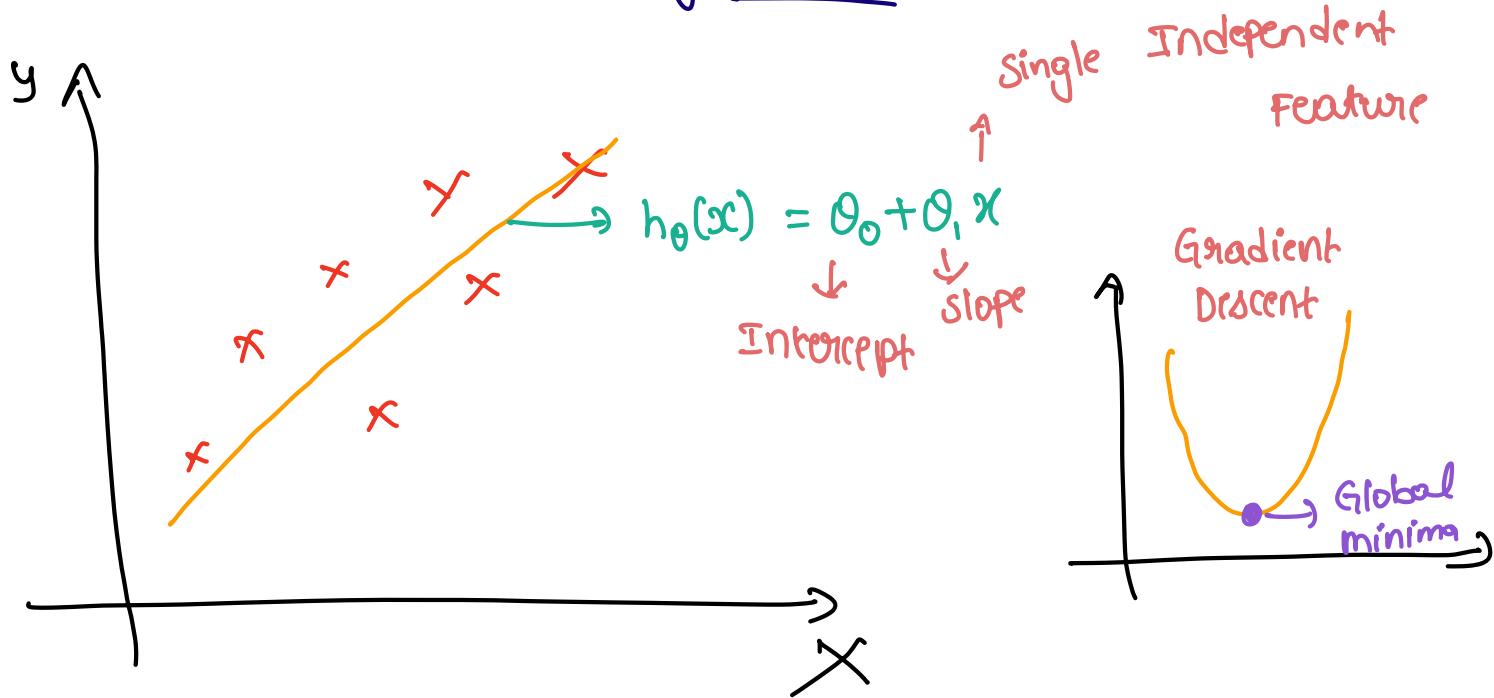
For Eg if we have 2 Independent Features

degree = 1 $h_{\theta}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

degree = 2, $h_{\theta}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2$

Since degree is 2, we have to square
Every input Feature

Summary of Linear Regression

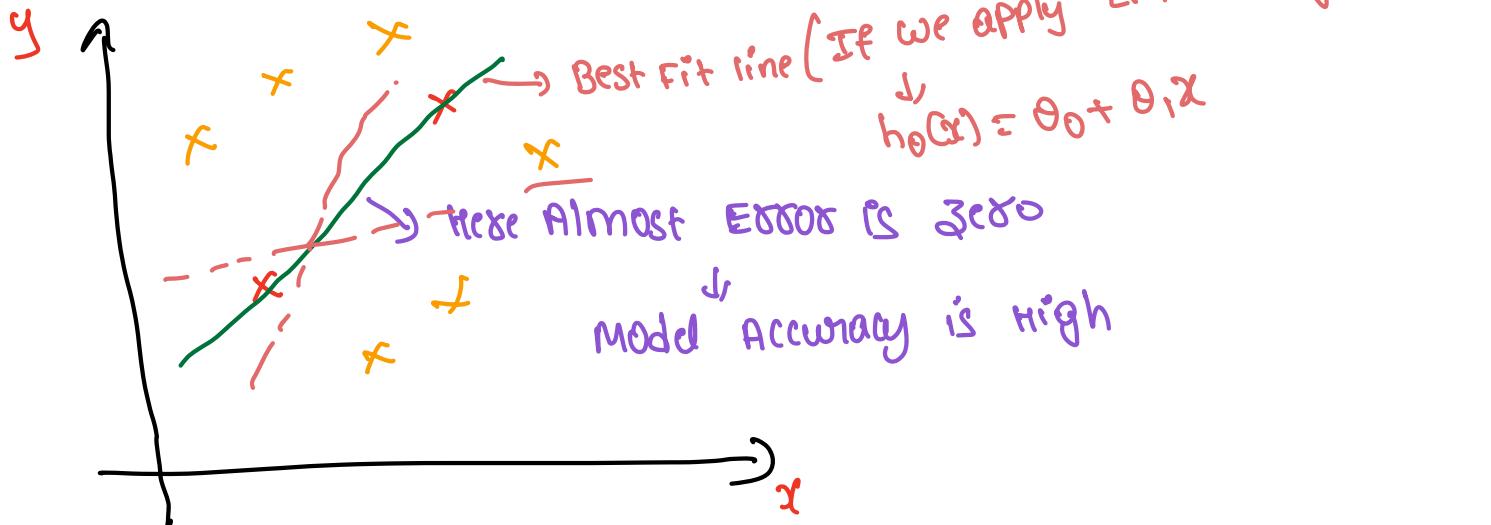


Cost function :

$$= \frac{1}{2m} \sum_{i=1}^m \left(\underbrace{h_{\theta}(x^{(i)})}_{\text{Predicted Value}} - \underbrace{y^{(i)}}_{\text{Actual Value}} \right)^2 \rightarrow \text{Mean square error}$$

Ridge Regression :

- ↳ Assume single Independent Feature
- ↳ 2 data points



$x \rightarrow$ New Test Data

↓
 Now the Error is Increasing on New Test Data

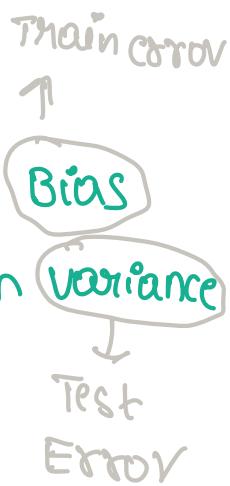
↓
 Model is Over Fitting

↓
 Train data → Accuracy ↑

Test data → Accuracy ↓ → High Variance

↓
 So, we use Ridge Regression

for Solving over fitting issue



① Ridge Regression (or) L2 Regularization :

↳ Reduces Overfitting

$$\text{Cost Function} = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

↓
 this Becomes
 zero during Overfitting

- So, we have to avoid the Cost Function shouldn't become zero
- So we will add 2 parts to the Cost Function as part of Ridge Regression

$$\text{Cost Function} = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{i=1}^n (\text{slope})^2$$

Hyper parameter

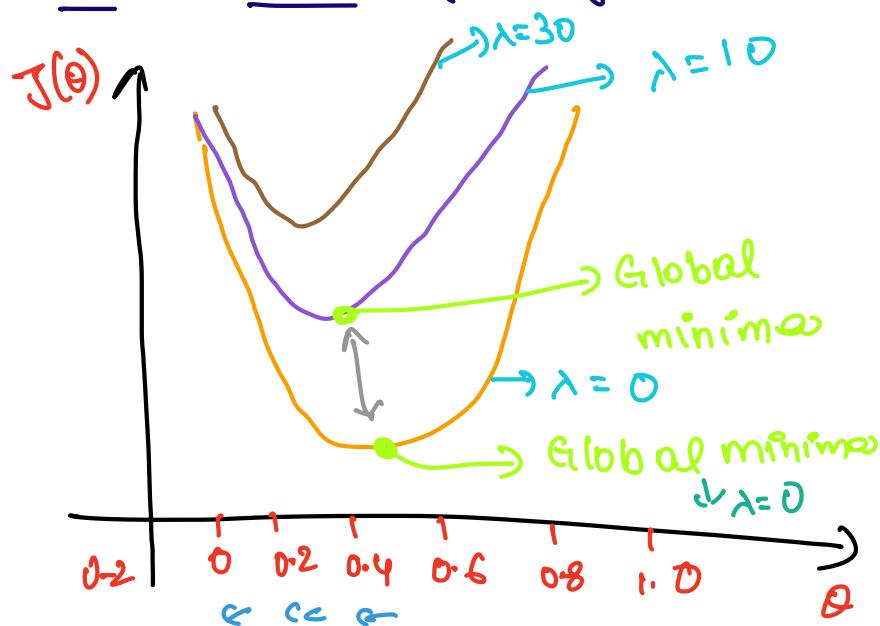
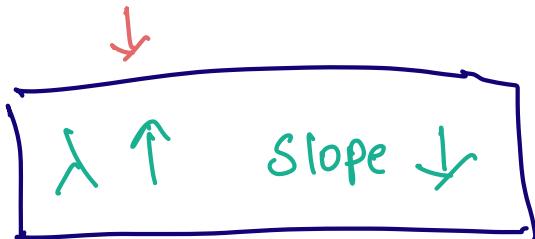
if $\lambda = 1$, $= (0) + (1) [(\theta_1)^2]$

↓
Here we are penalizing this '0' value by adding $(1) [(\theta_1)^2]$

↓
So the cost function will not be '0' at any time

↓
This will avoid over fitting

What is the Relationship B/w λ and slope² (L2 Reg term)



$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

Intercept ← = $0.34 + 0.52x_1 + 0.48x_2 + 0.24x_3$ (Assume)

If we Apply Ridge Regression on this

↓
Above Coefficients will reduce

↓
But will never become zero

$0.52x_1$

↓
with the unit movement in x_1 , what is the movement in y_1 ?

$$x_1 \rightarrow 1$$

$$y \rightarrow 0.52$$

↓
If x_1 is moving by '1'

↓
y will move by 0.52

For Eg :

In this case $x \rightarrow 1$ → If 'x' is moved by '1'
 $y \rightarrow 1$ → 'y' is moved by '1'
The coefficient is '1'

meaning 'x' and 'y' Features are highly correlated

After Applying Ridge Regression on Above :

$$h_0(x) = 0.34 + 0.40x_1 + 0.38x_2 + 0.14x_3$$

These coefficients are penalized (λ) rendered

→ How Ridge Regression (α) L2 reducing overfitting :



Reduces the impact of the coefficient, by reducing the coefficient of the feature that are not directly related, correlated with the output Feature

② Lasso Regression (L1 Regularization) → Feature selection

$$\text{cost function} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) + \lambda \sum_{i=1}^n |\text{slope}|$$

Adding this part to cost function as part of Lasso Regression

Relationship B/w ' λ ' and 'slope' :

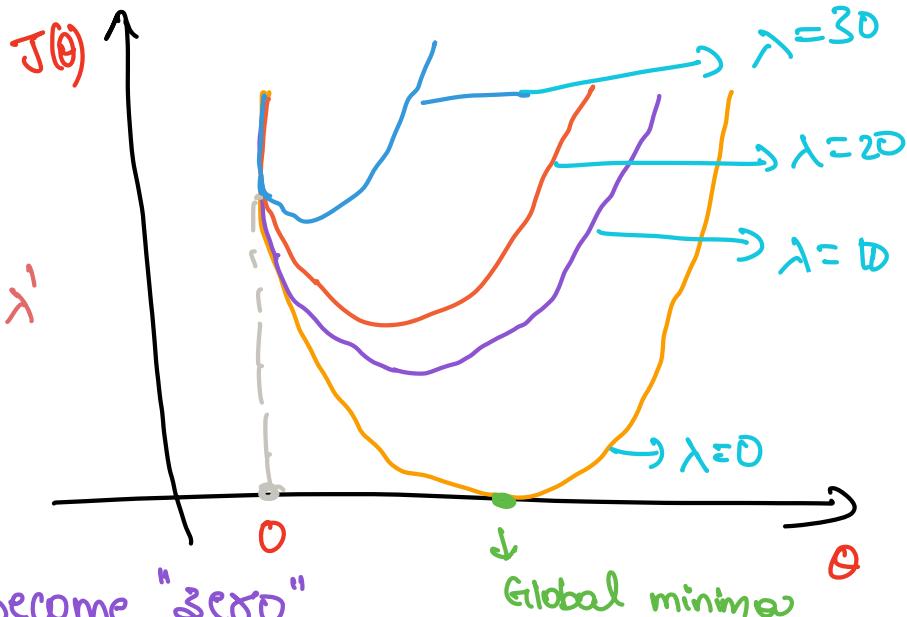
→ Here

$\lambda \uparrow, \theta \downarrow$

But at some point of ' λ '

$$\boxed{\theta = 0}$$

i.e., coefficient will become "zero"



i.e., In short we are trying to remove that specific Feature

For Eg :

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$
$$= 0.52 + 0.65x_1 + 0.72x_2 + 0.34x_3 + 0.12x_4$$

x_4 is the feature that is
not much correlated with O/p Feature

↓ ↓
Very Small coefficient 0.12

After Applying Lasso Regression (or) L1 on this
Cost Function

↓
Lasso will reduce it '0'

↓
'0.12' Becomes '0'

↓
 x_4 Feature is Eliminated

$$h_0(x) = 0.52 + 0.5x_1 + 0.60x_2 + 0.14x_3 + \cancel{0.12x_4}$$

→ Which all Input Features that are not much correlated
with Output Feature will have small coefficient value

↓
On Applying Lasso Regression, it reduces those
unimportant Features to 'zero'

↓
So we can use Lasso Regression for Feature selection

↓
If we have more Features, we can use Lasso

③ Elastic Net Regression

→ It is the combination of Both Ridge & Lasso

↳ ① Reduce overfitting

↳ ② Feature selection

$$\text{Cost Function} = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda_1 \sum_{i=1}^m (\text{slope})^2 + \lambda_2 \sum_{i=1}^n |\text{slope}|$$

\downarrow
L2
 \downarrow
Ridge Regression
 \downarrow
Reduce overfitting

\downarrow
L1
 \downarrow
Lasso Regression
 \downarrow
Feature Selection

→ For Hyperparameter Tuning

↓
We will use Ridge Regression, Lasso Regression
& Elastic Net Regression