

## Lecture - 01

→ Introduction to ML (AI vs ML vs DL vs DS)

→ Supervised ML and Unsupervised ML.

→ Linear Regression (Maths and Geometric Intuitions)

→  $R^2$  and Adjusted  $R^2$ .

→ Ridge and Lasso Regression.

### Introduction to ML:

AI vs ML vs DL vs DS

### Artificial Intelligence:

It is the simulation of human intelligence processes

by Machines, especially Computer Systems. Specific applications

of AI include Expert System, NLP, Speech Recognition and

Computer Vision.

### Machine Learning:

It is the branch of artificial intelligence (AI)

and Computer Science which focuses on the use of data

and algorithm to imitate the way human learn to

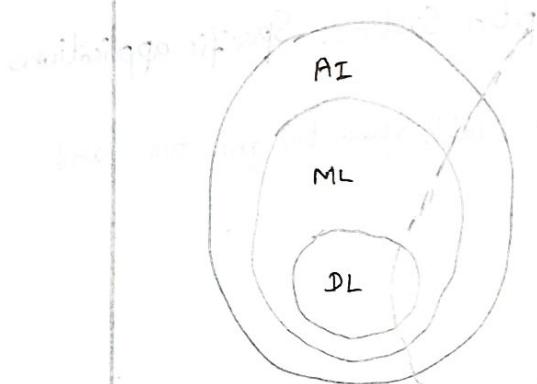
Machines.

## Deep Learning:

Deep learning is a method in artificial intelligence (AI) that teaches computers to process data in a way that is inspired by human brains. [Using Neural network]

## Data Science:

Data Science is the study of data to extract meaningful insights for business. It is multidisciplinary approach that combines principles and practices from the fields of Mathematics, statistics, artificial intelligence and Computer Engineering to analyze large amount of data.

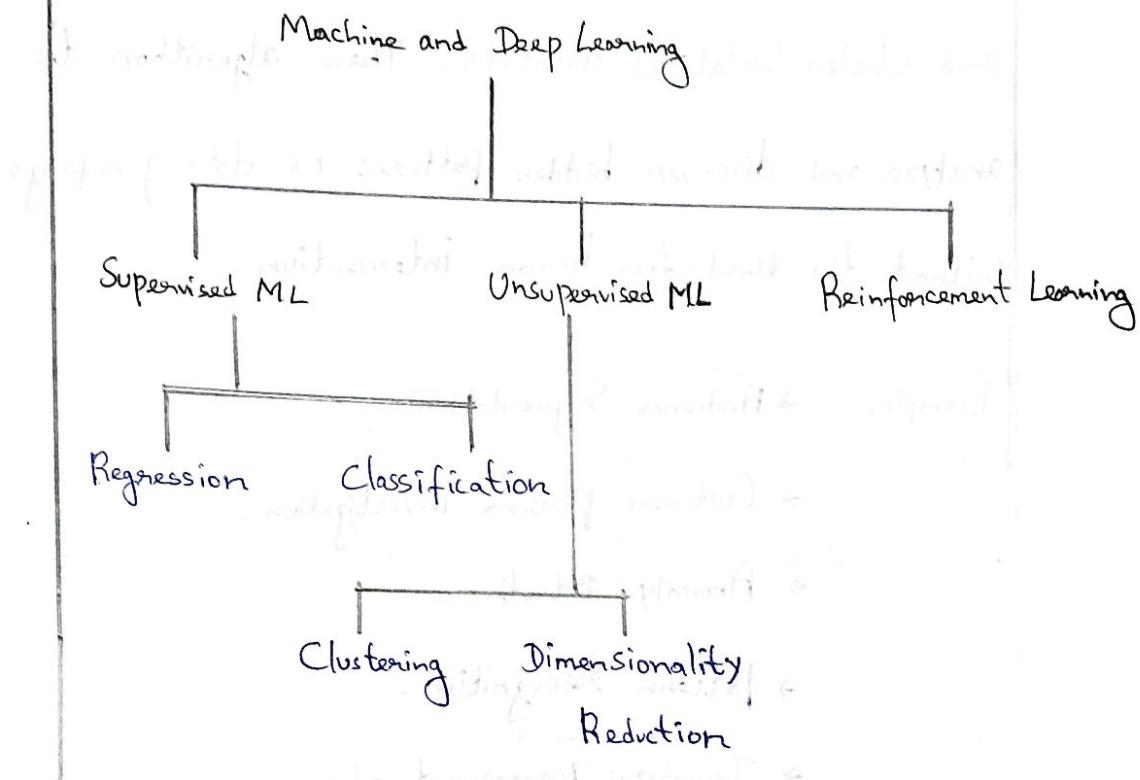


## (2) Examples,

Self driving Cars, Predictive modelling, Analyzing, Netflix Movie Recommendations etc.

## Supervised and Unsupervised Learning

about working of a general solution and AI



### Supervised Learning:

It is the Subcategory of Machine Learning and artificial Intelligence. It is defined by its Use of labeled datasets to train algorithms that to classify data (or) predict outcomes accurately.

Example, → Text Categorization .

→ Face Detection .

→ Spam detection .

→ Weather forecasting etc.

## Unsupervised Learning:

It uses machine learning algorithms to analyze and cluster unlabeled datasets. These algorithms to analyze and discover hidden patterns or data groupings without the need for human intervention.

Example, → Audience Segmentation.

→ Customer persona investigation.

→ Anomaly Detection.

→ Pattern recognition.

→ Inventory Management etc.

## Reinforcement Learning:

It is a machine learning training method based on rewarding desired behaviours and/or punishing undesired ones. In general, a reinforcement learning agent is able to perceive and interpret its environment, take actions and learn through trial and error.

Example,

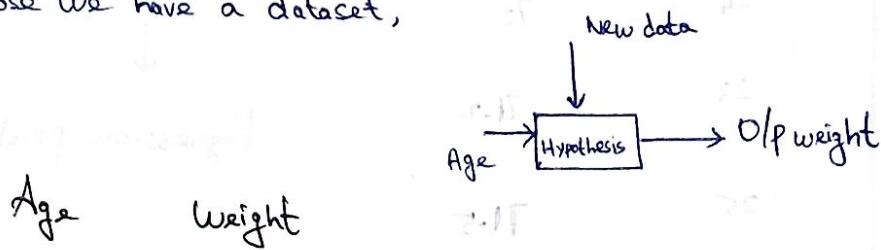
→ Playing a game [Chess, Go]

→ Humanoid Robots.

→ Self driving Cars etc.

## Supervised ML:

Suppose we have a dataset,



Age	Weight	We want to predict the weight respect to age. Suppose New input of age, we want to predict weight.
24	62	
25	63	
21	72	
27	62	

Independent Variable  $\rightarrow$  Age and other features in the dataset.

Dependent Variable  $\rightarrow$  Weight [Predicted feature]

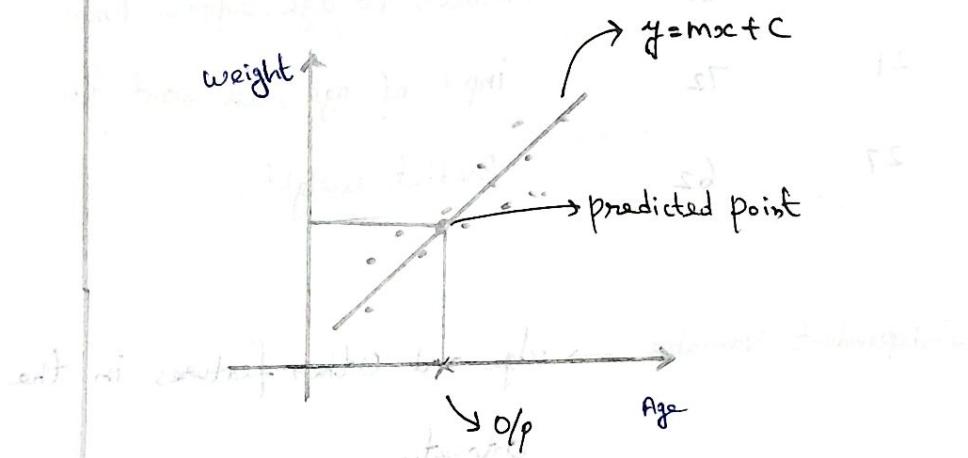
## Regression Problem:

It is the type of problem in Supervised learning, this type of problem cannot classify. Instead, this problem output as a continuous variable. These type of numerical output problem is called as regression.

Example,

Age	Weight	O/P
24	72	Continuous Variable
23	71.2	Regression problem
25	71.5	

If we are plot this, look like



Classification Problem:

In this problem, we can classify the output.

The Output dependent Variable is based on all or some independent Variable in the dataset. The classification

Problem is for Only Categorical Variable Output

Example,

No of hours study	No of hours play	No of hour sleep	Pass (Fail)
6	4	3	P
5	2	8	F
-	-	-	P
-	-	-	F

Final Decision Rule!!

} Binary  
Classification

↳ Output  
Categorical variable.

### Unsupervised ML:

The Unsupervised ML has no labeled one.

Suppose we have a dataset,

Person	Gender	Height
AA	Male	72
BB	female	72.5
CC	Male	73

The above dataset, There are only independent Variable.

There is no dependent Variable (or) any labeled feature.

So, this non-labeled one and this is called as the Unsupervised ML.

It is mainly focus on categorizing of data and  
Predict nearest point.

There are two category are there,

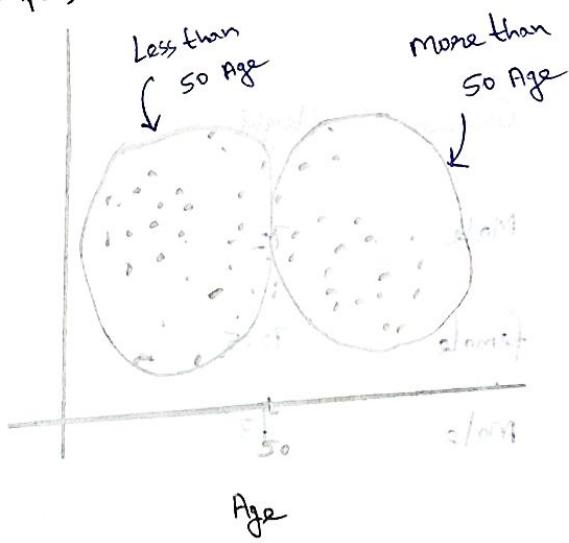
→ Clustering.

→ Dimensionality Reduction.

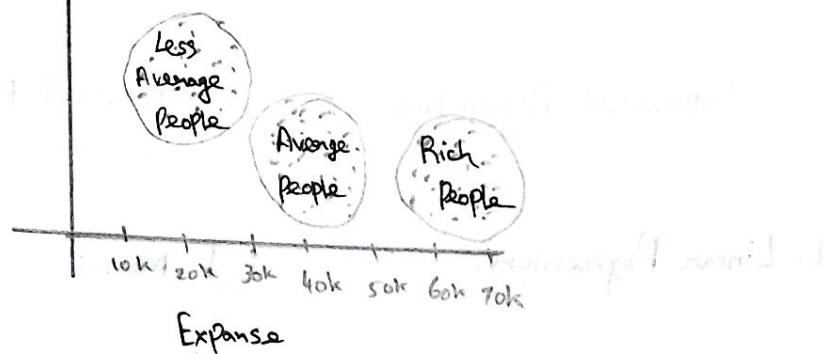
Clustering:

→ Clustering is nothing but grouping of the  
particular features on the dataset.

Example,



Customer Segmentation is one of the perfect example  
for clustering.



Based on this, The rich people have more expence in products. So again analyze products which products are highly purchased from rich people. Based on this, The Such Product Productivity is high.

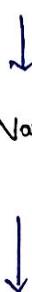
### Dimensionality Reduction:

The dimensionality reduction is nothing but the feature reduction.

Suppose we have 500 features, we reduce into 50 (or) less features.

#### Feature Reduction on dimension

Of dataset. (100 features)



Using Various algorithm

[PCA, LDA] etc.



Reduce features

## Discussing Algorithms:

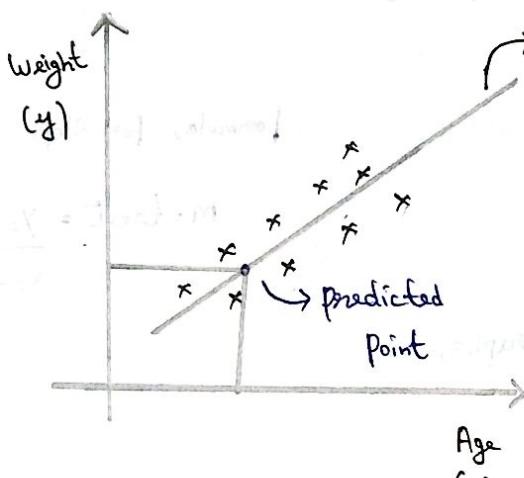
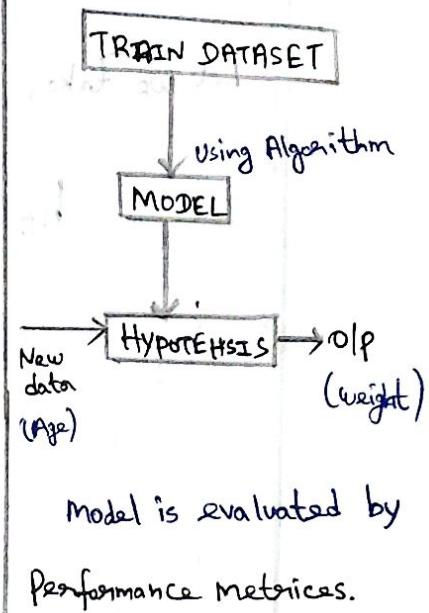
Supervised Algorithms	Unsupervised Algorithms
1. Linear Regression.	1. k Means.
2. Ridge and Lasso Regression.	2. DB Scan.
3. Logistic Regression.	3. Hierarchical Clustering.
4. Decision Tree.	4. k Nearest Neighbour Clustering.
5. Ada Boost.	5. PCA.
6. Random forest.	6. LDA.
7. Gradient Boosting.	
8. Xg boost.	
9. Naive Bayes	
10. SVM, KNN	

The above algorithm are the most important algorithm in machine learning and use it various applications.

## Linear Regression:

Linear Regression analysis is used to predict the value of a variable based on the value of another variable using constructing the straight line.

In normally,



$y$  is the linear function of  $x$ .

Equation of straight line,

$$y = Mx + C$$

Also written as,

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

(or)

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

## Equation of straight line:

$$y = M^{(i)} x + C$$

$$y = \beta_0 + \beta_1 x^{(i)}$$

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

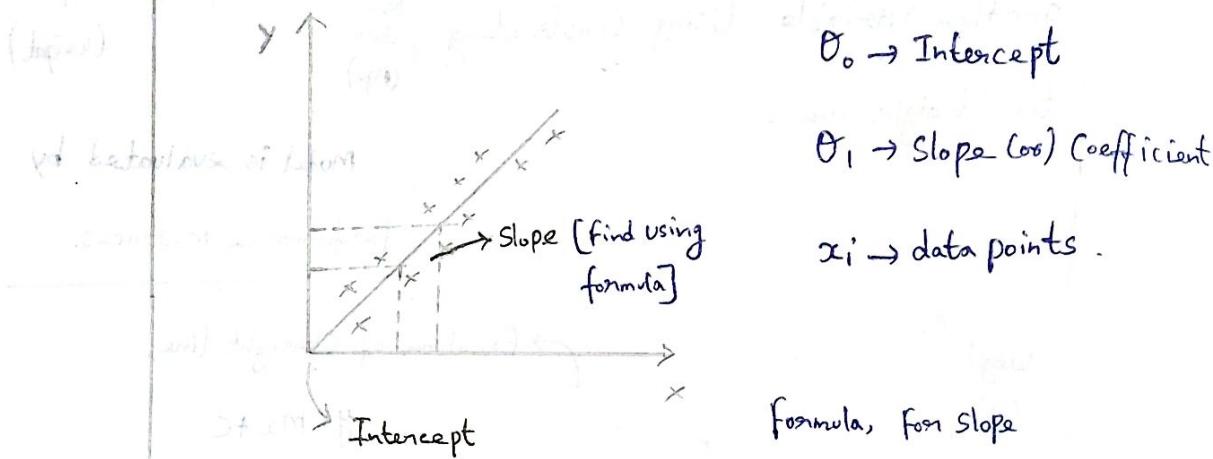
(i) → Indicates the independent Variable.

## Equation of straight line:

Let we take,

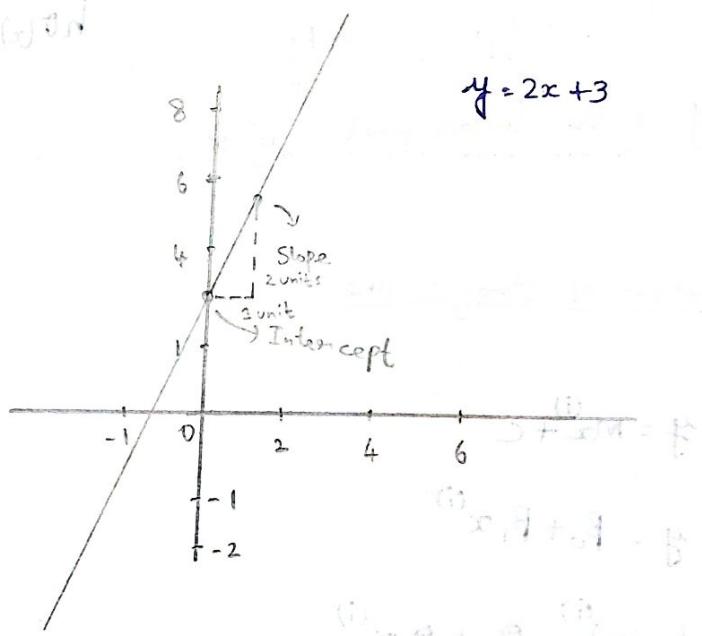
$$y_0(x) = \theta_0 + \theta_1 x$$

for value of  $y_0$  when  $x=0$



Equation for lines on graphs,

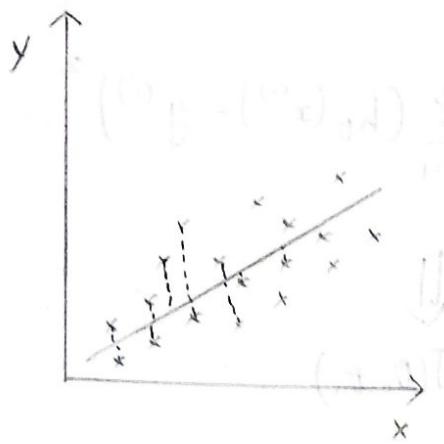
$$x_1 \theta + \theta = 0$$



The goal of linear regression is to find best fit

line across the datapoints.

How to find best fit line?

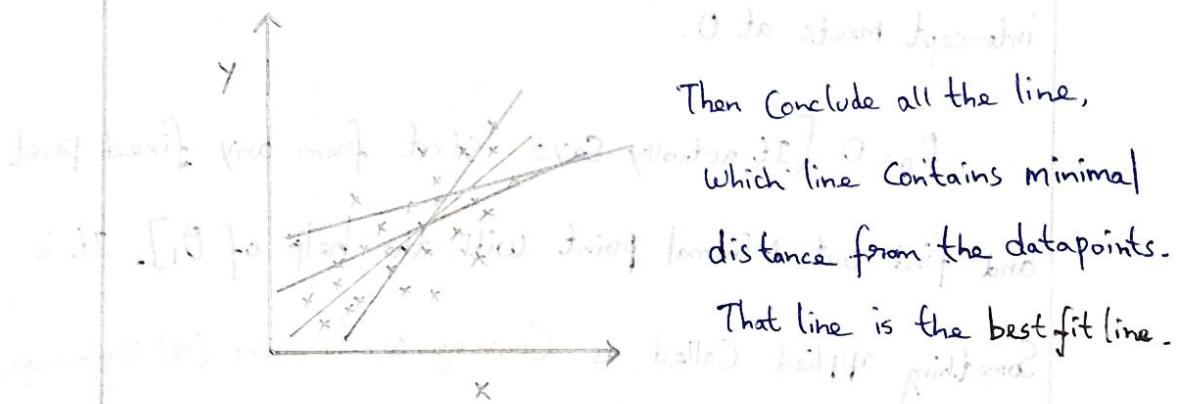


The datapoints from the line  
(distance) Should be Minimal.

The calculation is done by  
Cost function.

Actually, the Calculation is done every datapoints and  
Summation itself.

Samewise, We can Create Multiple linear line.



Cost function: [Distance formula]

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

Cost function,

$h_{\theta}(x) \rightarrow$  line Equation.

$y \rightarrow$  data points (Actual)

$\bar{x} \rightarrow$  Average of all Points.

$\sum \rightarrow$  Summation of all points.

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

Cost function      ↓  
Entire Equation called  
as Squared Error Function

$( )^2 \rightarrow$  Avoid neg value  
 $\frac{1}{2} \rightarrow$  For Simple derivative  
purpose.  
Suppose,

$$\frac{d(x^2)}{dx} = 2x \Rightarrow \frac{2x \times 1}{2}$$

What we need to solve?

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

$\theta_0, \theta_1$



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

We need to minimize [Cost function]

Now we calculate gradient with respect to  $\theta_0$  and  $\theta_1$ .

Consider Some Cases:

Before that take  $\theta_0 = 0$  for linear regression because the intercept meets at 0.

$\theta_0 = 0$  [It actually says start from any fixed point

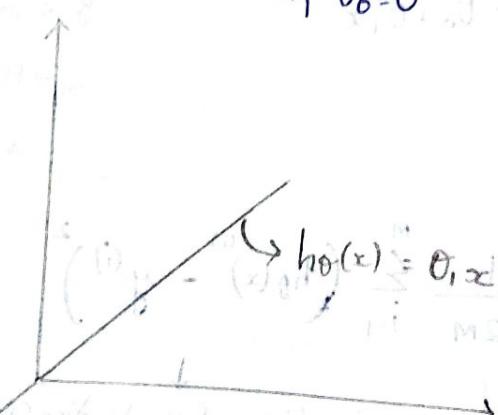
and find out minimal point with the help of  $\theta_1$ ]. It is

Something applied Called as Convergence theorem (or) Algorithm.

Case 1:

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

If  $\theta_0 = 0$



After fixing  $\theta_0 = 0$  and keep changing  $\theta_1$  value to find out best fit line with the help of gradient descent.

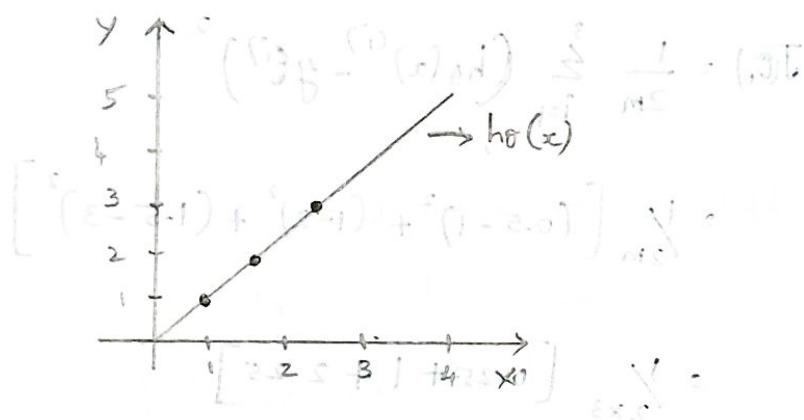
Case 2:

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

Suppose, we have the datapoints  $(1, 1), (2, 2), (3, 3)$ .

Consider the slope value, [case 2-1]

$$\theta_1 = 1$$



Try to find Minimal with respect to Cost function,

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^3 (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Apply datapoints  
Value.

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

$$= \frac{1}{6} [0 + 0 + 0] \Rightarrow 0$$

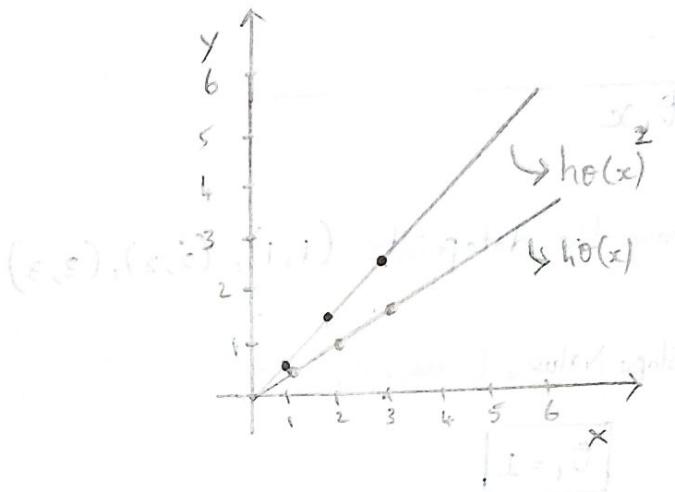
$$J(\theta_1) = 0$$

$$\theta_0 = (\theta_0)_0 \text{ Jd}, \theta_1 = (\theta_1)_0 \text{ Jd}, \theta_2 = (\theta_2)_0 \text{ Jd}$$

[Case 2.2] Find out the cost function

Consider the slope value,  $\theta_1 = 0.5$

$$\theta_1 = 0.5$$



$$h_{\theta}(x^{(1)}) = 0.5 * 1$$

$$h_{\theta}(x^{(1)}) = 0.5$$

$$h_{\theta}(x^{(2)}) = 0.5 * 2$$

$$h_{\theta}(x^{(2)}) = 1$$

$$h_{\theta}(x^{(3)}) = 0.5 * 3$$

$$h_{\theta}(x^{(3)}) = 1.5$$

Let's find  $J(\theta_1)$ ,

$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^3 (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \left[ (0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2 \right] \end{aligned}$$

$$= \frac{1}{2 \times 3} [0.25 + 1 + 2.25]$$

$$= 0.58$$

$$J(\theta_1) = 0.58$$

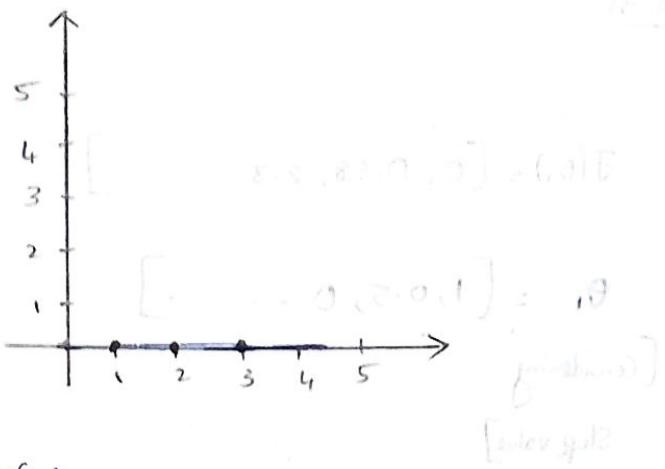
[Case 2.3]

Consider the slope value,  $\theta_1 = 0$

So,

$$\theta = (0, 0)$$

$$h_{\theta}(x^{(1)}) = 0, h_{\theta}(x^{(2)}) = 0, h_{\theta}(x^{(3)}) = 0$$



Let's find  $J(\theta_1)$ ,

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

$$= \frac{1}{6} [1+4+9] = 2.3$$

$$\boxed{J(\theta_1) = 2.3}$$

To try this number of times and get different different  $J(\theta_1)$  values. Based on this value, to find which one is minimal distance with the help of Gradient descent.

Real time:

We are just considering  $\theta(0)$  as fixed and try different  $J(\theta_1)$  value. The actual datapoints are also vary based on  $h_\theta(x)$ . [It means Actual datapoints  $x$  and  $y$  vary].

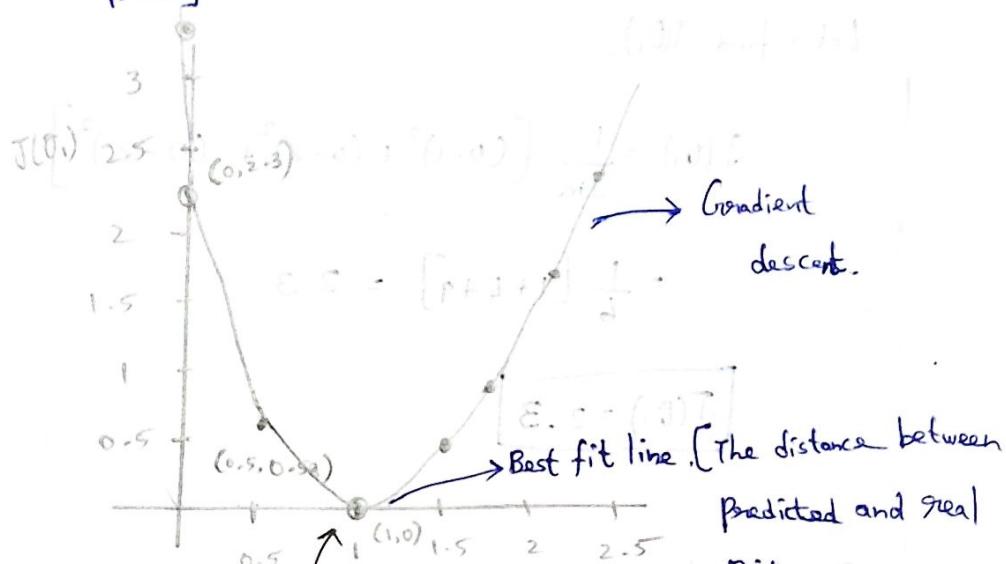
In real time, we just  $\theta(0)$  as fixed based on the graph and try different  $J(\theta_1)$  value. Finally find the minimal one Using Gradient descent.

Case 3:

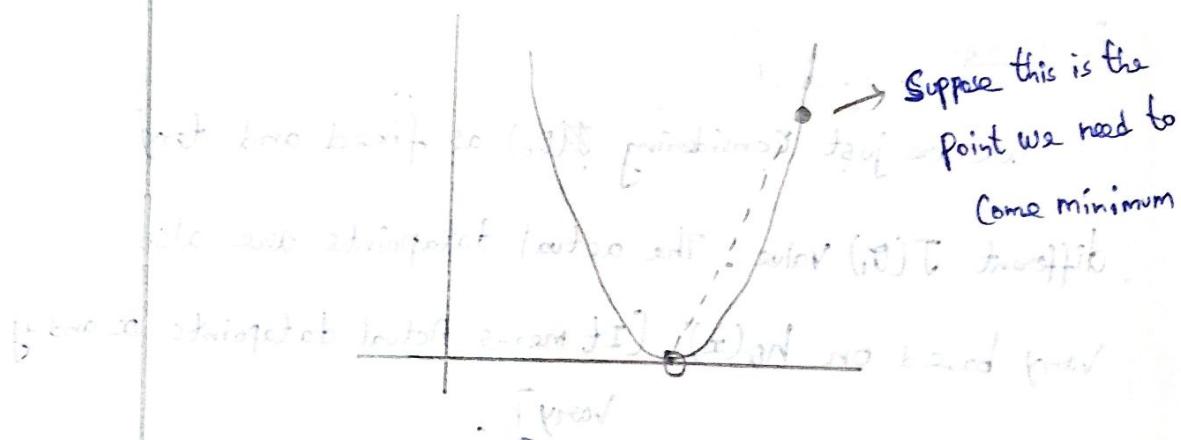
$$J(\theta_1) = [0, 0.58, 2.3 \dots \dots]$$

$$\theta_1 = [1, 0.5, 0 \dots \dots]$$

[Considering  
slope value]



J(theta\_1) vs theta\_1  
Best fit line [The distance between  
Predicted and real  
points are very very  
less]  
Global Minima



It is local cost function, and local cost function is not diff. with respect to theta\_1. It is local cost function, and local cost function is not diff. with respect to theta\_1. It is local cost function, and local cost function is not diff. with respect to theta\_1.

## Convergence Algorithm:

It is the algorithm for gradient descent. We can say [Repeat Until Convergence]

Algorithm,

Initial value of  $\theta$  to  $J^2$  taking value of  $\theta$  to  $J^2$

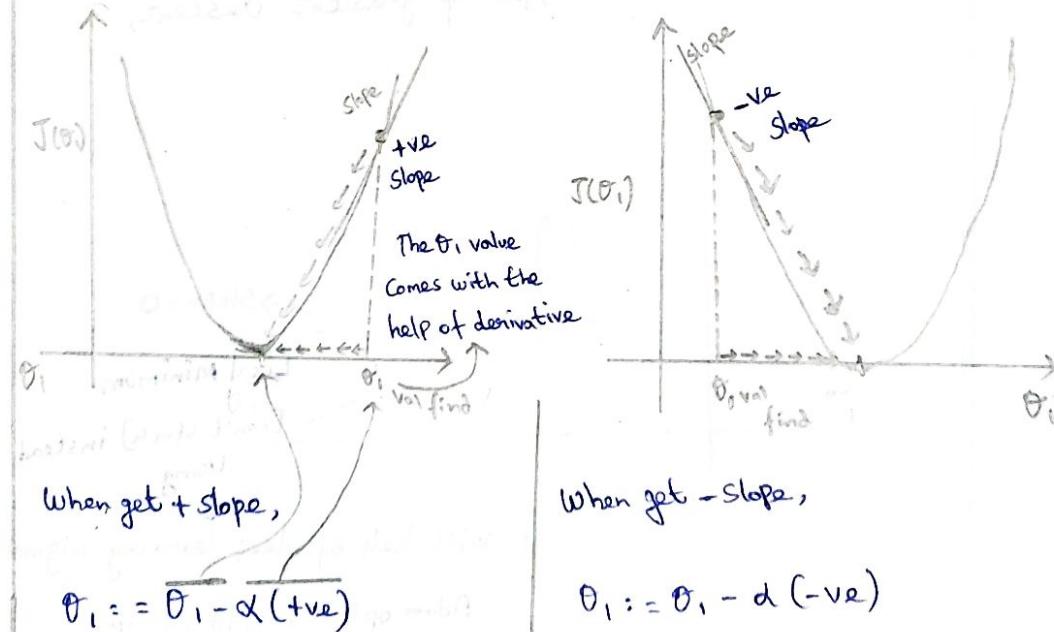
with Repeat Until Convergence end of not found

$$\left\{ \begin{array}{l} \text{initial value of } \theta \\ \text{repeat until convergence} \end{array} \right. \quad \begin{array}{l} \rightarrow \text{Derivative of} \\ \text{slope } \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \end{array}$$

$$\theta_j := \theta_j - \alpha \left[ \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \right] \quad \text{[Ans]}$$

3

Suppose we have two Cost function,



↳ This formula has various summation and iteration to reach global minimum.

$\alpha$  → Learning rate

When get -slope,

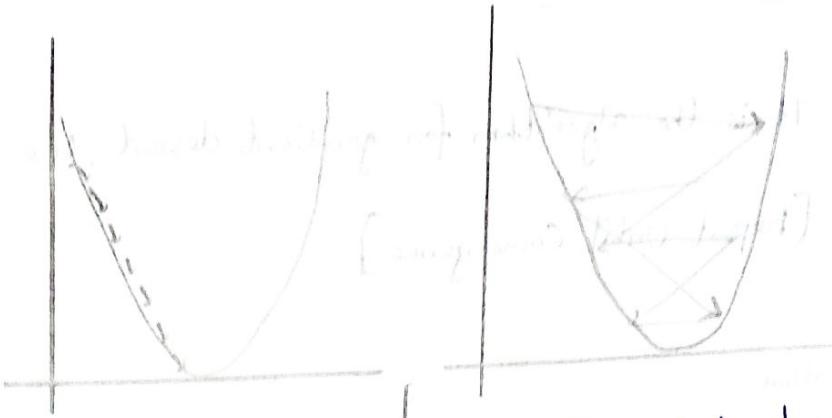
$$\theta_1 := \theta_1 - \alpha (-ve)$$

$$\theta_1 := \theta_1 + \alpha (ve)$$

$\alpha$  → Learning rate.

much How Speed the Solution to

find Global Minima is called as Learning rate.



If  $\alpha$  is too small, gradient descent can be slow.

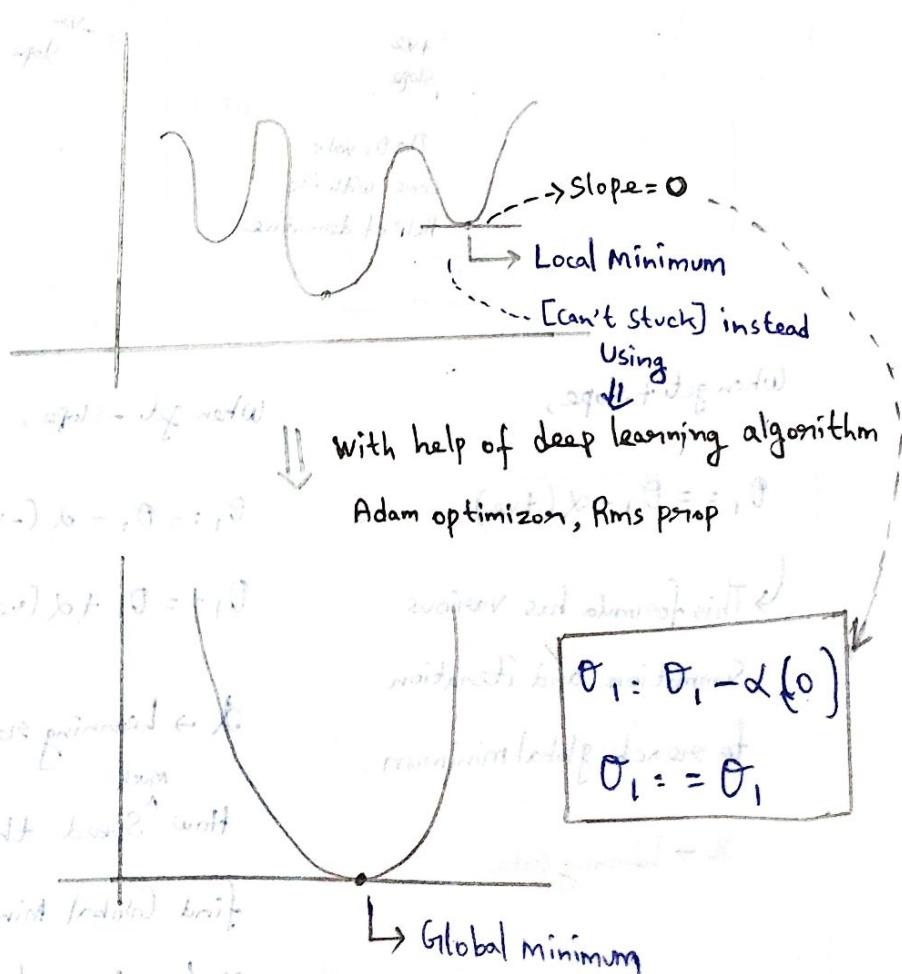
$$\alpha = 0.01$$

[But it takes time to reach]

If  $\alpha$  is too large, gradient descent can overshoot the minimum, It may to converge or even diverge.

Maintain Middle  $\alpha \rightarrow$  Learning rate is better.

Suppose, we have this type of gradient descent,



### Interview Question:

Does any local minima is present in linear regression?

In Machine learning problems Moreover see the Global minimum. In Some Cases, we are seen in local Minimum are found long convergence. At that time, Using deep learning algorithm (or) optimizers to solve local Minimum to reach Global Minimum.

Convergence algorithm is nothing but a Gradient descent.

$$\text{Algorithm, } \theta^{(t+1)} = \theta^{(t)} - \alpha \sum_{i=1}^m [y_i - h_{\theta}(x_i)]$$

Repeat until Convergence

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

Find Out the  $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ ,

When  $j=0$  and  $1$  [Because  $\theta_0$  and  $\theta_1$ ] i.e  $h_{\theta}(x) = \theta_0 + \theta_1 x$

Rewrite us,

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

Partial derivative applied

$j=0$ ,

↓

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

$\frac{1}{m} x^2$

$$\begin{aligned} & 2 \times \frac{1}{m} \\ & = \frac{1}{m} \end{aligned}$$

(i)  $\theta_0$  movement calculate forward without  $\theta_1$

$j=1$ ,

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

$$h_\theta(x) = \theta_0 + \theta_1 x_i$$

Repeat until Convergence

{

backward and parallel to nothing and nothing

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

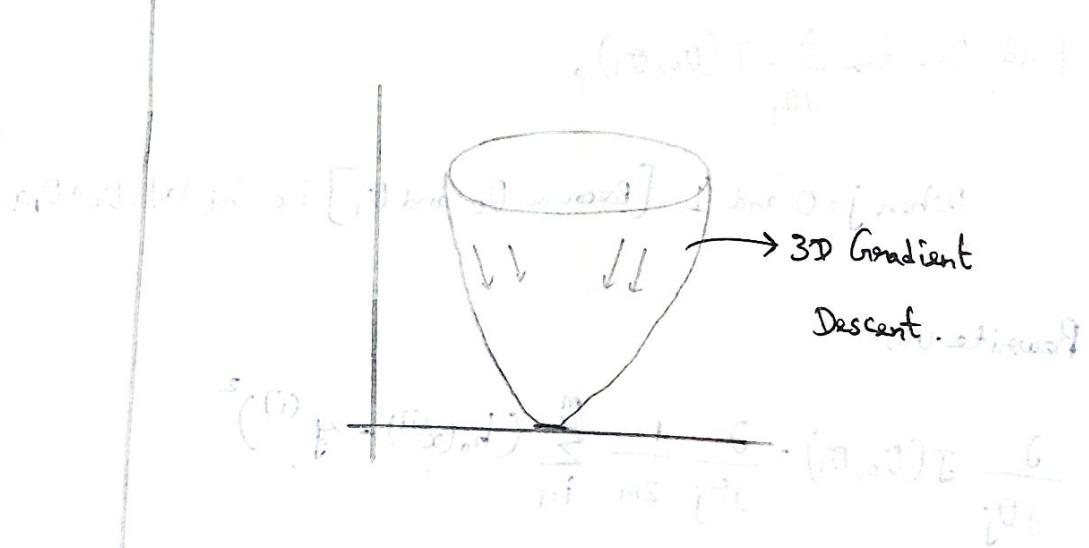
forward then repeat

This is known as the linear regression.

$$(3, 3) \rightarrow \frac{6}{3} = 2 - 2 = 0$$

Suppose we have multiple number of features  $x_1, x_2, x_3, x_4, \dots$

$x_n$ . The Gradient descent look like,



## Performance Metrics

$R^2$  and Adjusted  $R^2$ .

Performance metrics says that the how much this model can be efficient.

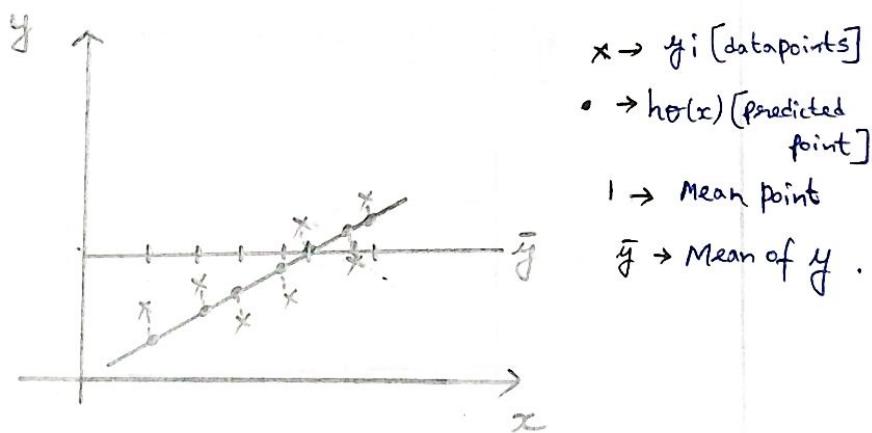
$R^2$ :

$$R^2 = 1 - \frac{SS_{\text{Residual}}}{SS_{\text{Total}}}$$

$\rightarrow$   $\hat{y}_i = h(x)$

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

$\rightarrow$   $\bar{y} \rightarrow$  Mean of  $y$ .



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

[This distance value is low]

[This distance value is high]

$$1 - \frac{\text{low}}{\text{high}}$$

$\rightarrow$  Small number with 0 to 1.

$$= 1 - \frac{\text{Low}}{\text{High}}$$

$$= 1 - \text{Small number} = 0.\text{Something} \quad [\text{Positive number}]$$

Big and small start with random numbers.

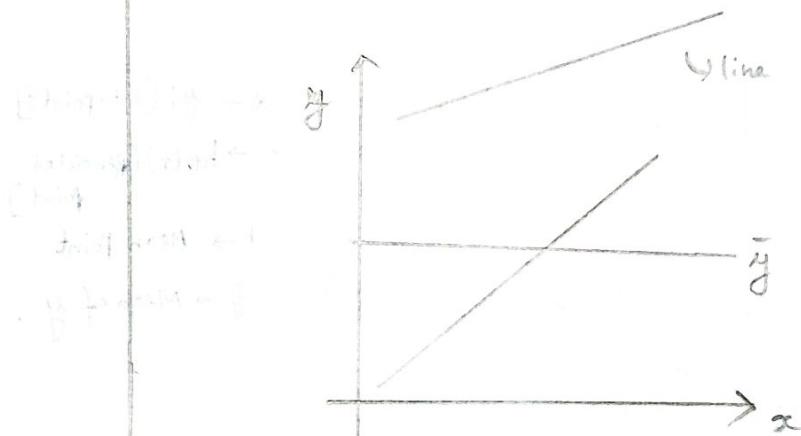
$$\text{Suppose, } 1 - \text{Small number} = 0.9$$

90% of accuracy.

$$\text{Suppose, } 1 - \text{big number} = \cancel{-\text{Ve number}}.$$

This Case is probably Occuring in Some Situations,

At that time. The graph look like,



In this case, Very less number of Chances are having.

Suppose we have dataset,

Dataset features:

Bedrooms      Price       $[R^2 = 85\%]$

Bedrooms      location      Price       $[R^2 = 90\% \uparrow]$  increases

because 'location'

Suppose we have one or two uncorrelated features, features correlated to price.

Gender      Bedrooms      Location      Price       $[R^2 = 91\%]$

$R^2$  increases, But Gender is not correlated to

Price

So,  $R^2$  take 91% as efficiency. But this is wrong efficiency.

Therefore, we move to the adjusted  $R^2$ .

Adjusted  $R^2$ :

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

P = features (or) predictors

N = No of datapoints.

if  $P=2$   $R^2 = 90\%$  but  $R^2 \text{ adjusted} = 86\%$ .  
↓ Increase  
↓ decrease

if  $P=3$   $R^2 = 91\%$  but  $R^2 \text{ adjusted} = 82\%$

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

↑ Big number  
↓ Small number

when  $P > 3$  (or) any feature the whole

$N-p-1$  keeps on decreasing.

So, The  $R^2$  is little bit increasing when the uncorrelated feature not huge. The adjusted  $R^2$

Solve that particular problem.