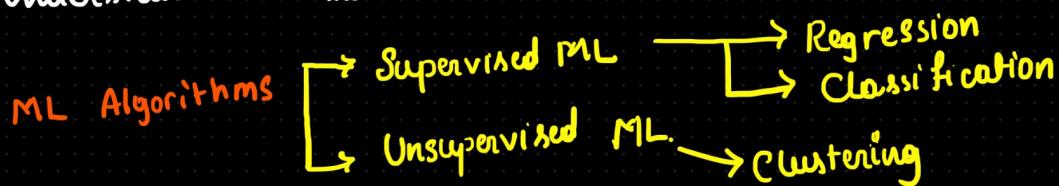


# Machine Learning Algorithms

In ML Intuition is very important for the interviews.

Understand the maths behind all algorithms.



Lets talk about Regression and classification problem

## Regression Problem

→ House price prediction.

No. of Rooms → Input ] → Independent variables

Total size → Input ]

Location → Input ]

price → Output. ] → Dependant variables

↳ 1.5m, 1.66m, 1.725m

↳ Output is continuous in regression problem

## Classification Problems

① Binary Classification

② Multiclass Classification

Example for binary classification

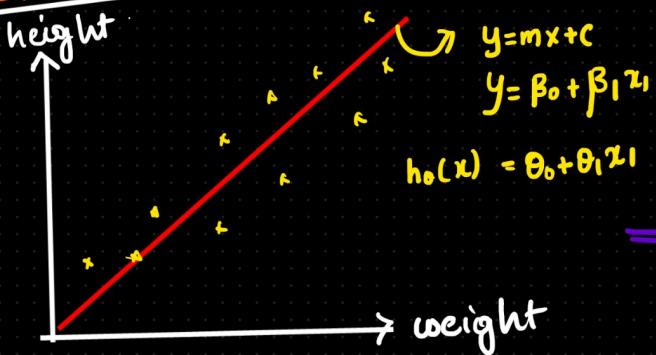
No. of play hour	No. of sleep hours	No. of study Hours	Pass/Fail	
5	8	2	F	→ Binary Classification
6	8	1	F	
0	7.5	5.5	P	
3	9	2	P	
4	5	7	P	

## Regression Algorithm:

1. linear Regression [Regression]
2. Ridge Regression [—, —]
3. Lasso Regression [—, —]
4. Elasticnet Regression [—, —]
5. logistic Regression [Classification]
6. Decision tree
7. Random forest
8. AdaBoost
9. XG Boost

Both Regression and Classification Problem.

## Simple Linear Regression

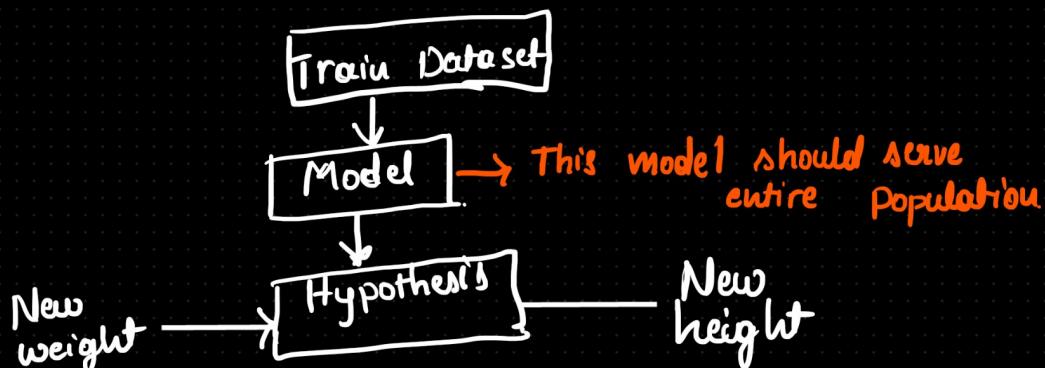


weight → I/P

height → O/P.

Goal of linear Regression is to create the best fit line.

⇒ If we have 2 feature we will call as simple linear Regression  
 ⇒ If we have more than 2 feature we will call as multiple linear Regression

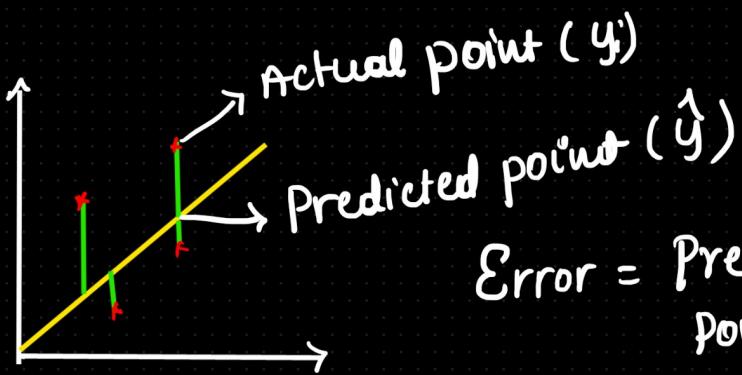


Equation of straight line

$$h_0(x) = \theta_0 + \theta_1 x \rightarrow \text{Equation used for hypothesis testing}$$

$\theta_0 \rightarrow$  Intercept

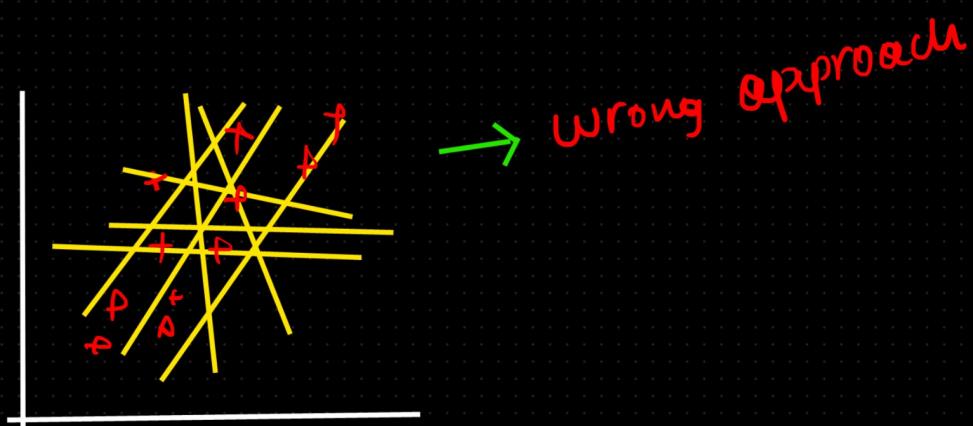
$\theta_1 \rightarrow$  Slope/Coefficient



$$\text{Error} = \text{Predicted Point} - \text{Actual Point}$$

$$\text{Error} = \hat{y} - y_i$$

- 1) we should change the value of  $\theta_0$  and  $\theta_1$  and we need to find best fit line with minimal error.



- 2) we cannot find best fit line by selecting random values. we need to choose a value and we need to go slowly towards best fit line.

- 3) To achieve best fit line we need to use the following 2 important concepts

⇒ COST FUNCTION

⇒ GRADIENT DESCENT

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 \Rightarrow \text{Now Error mean } \underline{\hat{y} - y}$$

$$= \text{COST FUNCTION} \quad J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$$

↓  
MEAN SQUARED ERROR

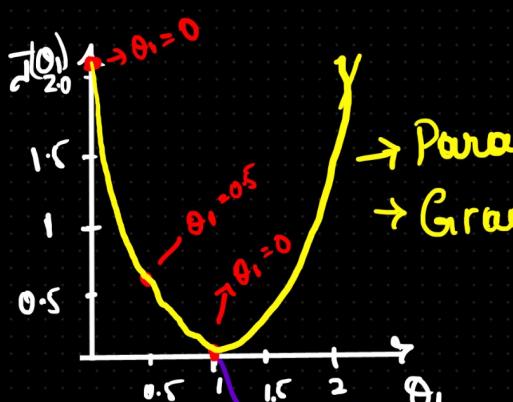
what we need to solve?

we need to make sure that

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

④  $h_\theta(x) = \theta_0 + \theta_1 x_1$  if my  $\theta_0 = 0$  my straight line

$$h_\theta(x) = \theta_1 x_1$$

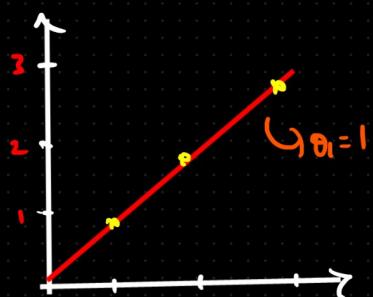


$$\Rightarrow \frac{1}{2m} \sum_{i=1}^m (y_i - y)^2$$

↓  
Quadratic Equation

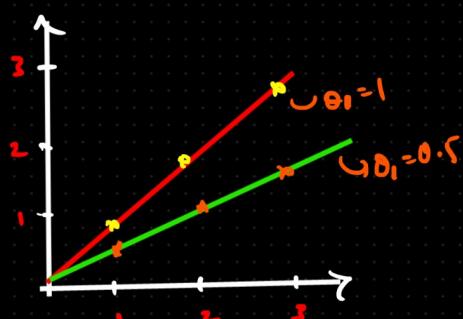
$$ax^2 + bx + c = 0$$

If  $\theta_1 = 1$

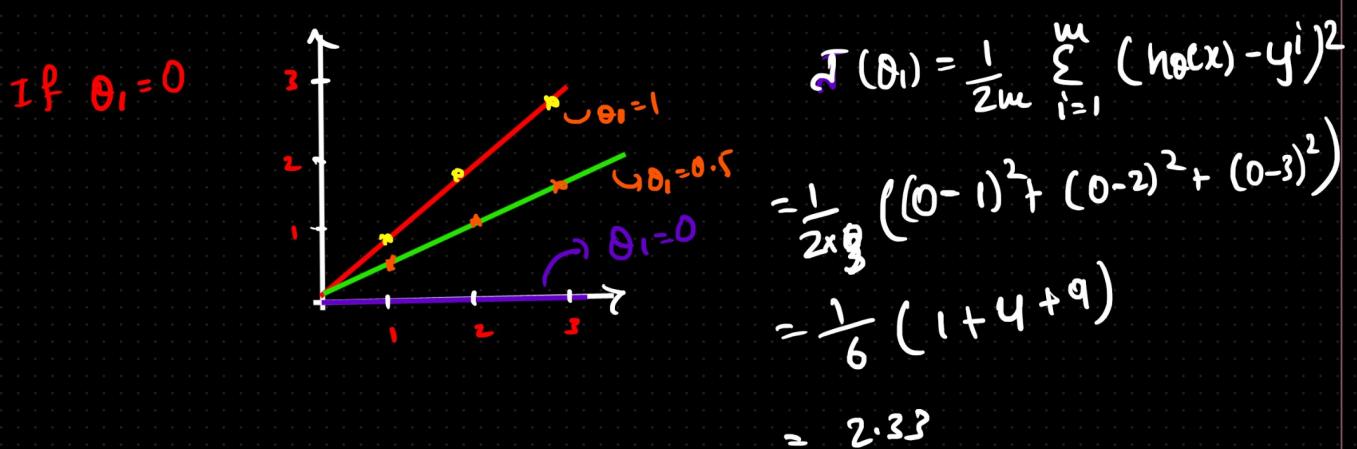


$$\begin{aligned} \bar{J}(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \\ &= \frac{1}{2m} \left[ (1-1)^2 + (2-2)^2 + (3-3)^2 \right] \\ &= 0 \end{aligned}$$

If  $\theta_1 = 0.5$



$$\begin{aligned} \bar{J}(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \\ &= \frac{1}{2m} \left[ (1-0.5)^2 + (2-1)^2 + (3-1.5)^2 \right] \\ &= \frac{1}{6} (0.25 + 1 + 2.25) \\ &= 0.58 \end{aligned}$$



## Gradient Descent (Convergence Algorithm)

repeat until

$$\left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial (J(\theta))}{\partial \theta_j} \\ \theta_{\text{new}} = \theta_{\text{old}} - \alpha \frac{\partial (J(\theta_{\text{old}}))}{\partial \theta_{\text{old}}} \end{array} \right.$$

$\theta_{\text{new}} \approx \theta_{\text{old}} - \alpha \text{ (-ve)}$

$\theta_{\text{new}} \approx \theta_{\text{old}} + \alpha \text{ (+ve)}$

$\theta_{\text{new}} \gg \theta_{\text{old}}$

**Learning Rate:** Learning used to determine the speed of convergence. Usually learning rate should be very small 0.001

In case of multi-linear Regression your formula will be

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \theta_5$$

### Advantage of MSE

- 1) It is having one global minima.
- 2) It is differentiable

### Disadvantage of MSE

- 1) Not robust to outliers
- 2) MSE Penalizes the error.

## Mean Absolute Error

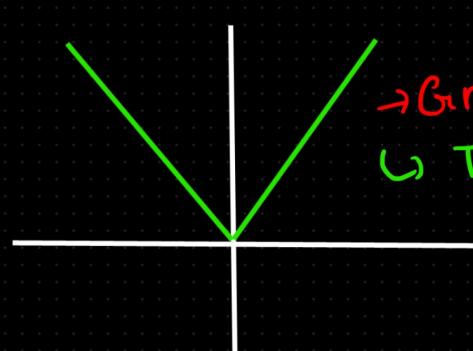
$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m |\hat{y}_i - y_i|$$

Note: If there is outlier then you can use MSE.

If there is outlier we have 2 option

→ Use MAE

→ Remove Outlier use MSE



→ Gradient Descent for MAE  
↳ This is called as subgradient

## Performance Matrix

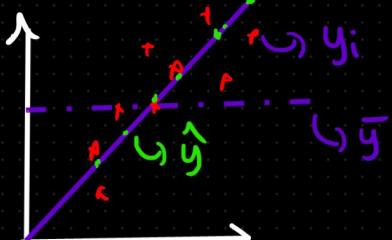
There are 2 performance matrix we are using

- 1)  $R^2$
- 2) Adjusted  $R^2$

$$R^2 = 1 - \frac{SS_{\text{Res}}}{SS_{\text{total}}}$$

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

$= 1 - \frac{\text{lower value}}{\text{higher value}}$  So always performance matrix  $R^2$  lies between 0 and 1.



## Important Interview Question:-

Can  $R^2$  be negative value?

Yes if can. and this means the best fit line is worst than average line

## Adjusted R<sup>2</sup>

House Pricing → Let's Consider

No. of Room    Total Size    Location    Gender    Price.



100

Correlated to Price      Not Correlated to Price

→ In  $R^2$  it will reduce the performance if there is any unrelevant feature are there

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

$N \rightarrow$  No. of data point

$$P=2 \approx R^2 = 90\%. R^2 \text{ Adj} = 86\%$$

$P \rightarrow$  No. of Predictors

$$P=3 \approx R^2 = 91\%. R^2 \text{ Adj} = 82\%$$

$(N-P-1)$  ↓  $\begin{cases} \rightarrow \text{If } P=2 \text{ denominator value will be higher} \\ \Rightarrow \text{If } P=3 \text{ denominator value will decrease} \end{cases}$   
denominator.

Let's take an example when  $P=2$

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

$$R^2 \text{ Adj} = 1 - \frac{(1-0.9)(100-1)}{100-2-1}$$

$$= 1 - \frac{0.1 \times 99}{97}$$

$$= 1 - \frac{9.9}{97}$$

$$= 1 - 0.102$$

$$= 0.8979$$

Let's take an example when  $P=3$

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

$$R^2 \text{ Adj} = 1 - \frac{(1-0.9)(100-1)}{100-3-1}$$

$$= 1 - \frac{0.1 \times 99}{96}$$

$$= 1 - \frac{9.9}{96}$$

$$= 1 - 0.1031$$

$$= 0.8969$$

So always  $R^2 > \text{Adjusted } R^2$

Conclusion:- we need to check always  $R^2$  and adjusted  $R^2$  both

If my difference between  $R^2$  and Adjusted  $R^2$  is higher then it means then there are many features that are highly correlated.

Linear Regression Assumptions:-

- ① There is a linear Relationship between  $x$  &  $y$
- ② Independent features should have Normal distribution
- ③ Always take care multi collinearity
- ④ Homoscedasticity = All the features have same variance
- ⑤ Feature Scaling required

RMSE (Root mean Squared Error) Gradient Descent

$$MSE = \frac{1}{2m} \sum_{i=1}^m \sqrt{(h_\theta(x) - y_i)^2}$$

$$RMSE = \sqrt{MSE}$$