

Linear Regression Algorithm

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

Multiple Linear Regression

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

Mean Square Error (MSE)

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

↑
Cost function

Q. What is the difference b/w Loss function and Cost function.

<u>Loss function</u>	<u>Cost function</u>
<p>① In case of loss function it will be w.r.t every observation.</p>	<p>① In cost function we will be doing calculations w.r.t all the data points. i.e. $i=1 \dots m$.</p>
<p>② In case of loss function we do the calculation w.r.t each & every point.</p>	<p>② Here we calculate the distance b/w predicted value & actual value for all the data points. basically we do the summation.</p>
<p>③ Loss function = $(h_{\theta}(x^{(i)}) - y^{(i)})^2$</p> $L.F = (\hat{y}_i - y_i)^2$ <p style="text-align: center;">↑ Predicted value Actual value</p>	<p>③ Cost function = $\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$</p> <p>④ Summation of square of difference of each and every point</p>
<p>④ Single Point difference.</p>	<p>⑤ Here after finding the difference and square of it we do the summation of them.</p>
<p>⑥ Calculate the difference of every point independently and then do square of it.</p>	<p>⑦ Here after doing finding the difference and square of it we do the summation of them.</p>

Convergence Algorithm

Repeat until Convergence

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

where $j = 1 \text{ to } m$

(θ_0, θ_1)

Find the Value of this derivative

$$\frac{\partial}{\partial j} J(\theta_j)$$

θ_0, θ_1

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

When,

$j=0$

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

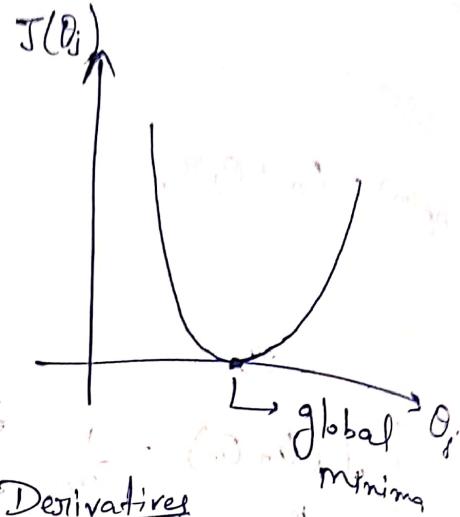
$$= \frac{\partial}{\partial \theta_0} \left[\frac{1}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})^2 \right]$$

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

when $j=1$

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

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



Derivatives

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

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

$$\begin{aligned} \frac{\partial}{\partial x} (x+1)^2 &= 2(x+1) \\ &= 2(n+1) \end{aligned}$$

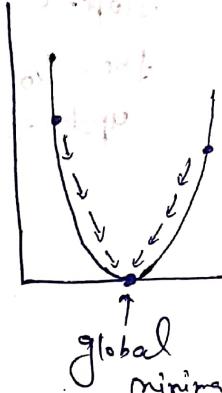
Now θ_0 and θ_1 Value will change according to this formula,

Repeat until Convergence

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

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

α determines the speed of convergence
and By Using this the loss will be minimum.



Types of Cost function

- i) MSE (Mean Squared Error)
- ii) MAE (Mean Absolute Error)
- iii) RMSE (Root Mean Square Error)

i) MSE (Mean Squared Error)

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

quadratic eqn.

Where $\hat{y} = \theta_0 + \theta_1 x$

\uparrow

Predicted value

We know that,

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

i.e it is

$$a^2 + b^2 + c = 0$$

→ This is a quadratic equation

→ And we plot the graph of a quadratic eqn it gives bell curve or convex function.

Advantage and Disadvantage of MSE

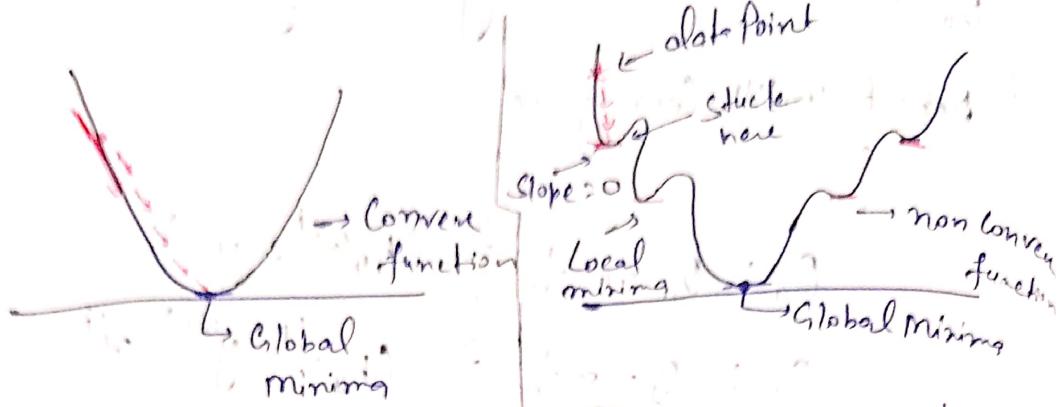
i) This eqn is differentiable.

Q Why do we need differentiability?

To calculate θ_0 , θ_1 and slope.

{differentiable means the derivative exists at every point in its domain. Consequently, the only way for a derivative to exist is if the function also exists (i.e continuous) on its domain. Thus a differentiable fn is also a continuous fn}

(ii) This eqn also has only one global minima.



* At the point of global minima $\text{slope} = 0$, and if slope is 0 then θ_0 and θ_1 will not get updated.

* Hence our aim is to always work with Convex function so that we can get only one global minima.

* Disadvantage of MSE:

i) It is not robust to outliers.

Q: Bez of outlier whether the cost fn will increase or decrease?

Soln: It will increase by tremendous amount (↑↑↑).

How?

* Bez of outlier the error will increase and θ_0 and θ_1 will get adjusted and since we are squaring the error the MSE will increase tremendously.

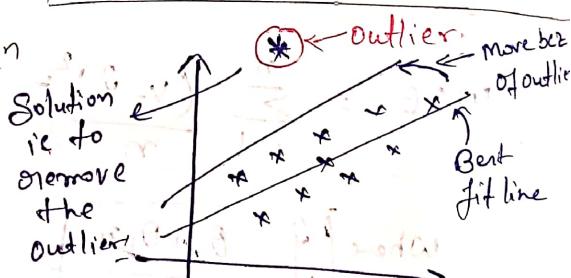
* Hence to solve this problem, we need to remove the outliers.

* Hence if we use MSE directly it will not robust to outliers because there will be bigger switch bez of outliers.

* In this case at local minima also $\text{slope} = 0$. Hence the θ_0 and θ_1 will not update further and it will get stuck at local minima.

* Hence this scenario will Meyer Come w.r.t. MSE

* Hence we can't use non-convex function for Cost function.



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

Experience (X)	Salary (Y)	Predicted Salary (\hat{y})	Error ($y - \hat{y}$)
1	100	110	-10
2	100	110	-10
3	100	110	-10
4	100	110	-10
5	100	110	-10

Suppose we are trying to predict the salary based on experience. In that case,

Let's calculate $(y - \hat{y})^2$

i) Here Error present in $(\text{lakh})^2$, i.e. $(\text{Salary} - \text{Predicted Salary})^2$ and bcz of that unit is changing and Error is getting Penalized

$$= (\text{lakh})^2$$

Eg: $(100 - 110)^2 = 10^2 = 100$

$\uparrow \quad \uparrow \quad \uparrow$

3digit 3digit 3digit

Hence Unit will change from lakh to millions.

ii) and Based on that we need to reduce the cost function.

but suppose $(999 - 100)^2, (899 - 100)^2, \dots, (100 - 100)^2, \dots, (100 - 100)^2 = 8021$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$

3digit 3digit 3digit 4digit

(unit changed)

* Bcz of changing of unit time complexity will increase and that's why outlier affect will also happen.

i) We are suppose to reduce MSE but here in this case it is increasing tremendously. Hence it is Penalizing the error i.e. changing the unit. Hence it is not robust to outliers.

ii) Penalizing the Error i.e. Changing the Unit.

Hence to solve this problem we ignore the outliers.

② Mean Absolute Error (MAE)

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

Advantage and Disadvantage of MAE

① Advantage of MAE

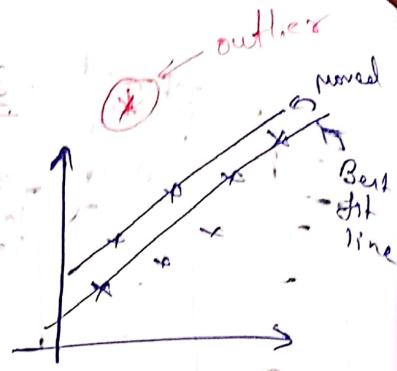
① Robust to outliers.

② It will also be in the same unit.

Since here we are not squaring the error.

Hence, bcz of outlier even if best fit line moves it will change little bit only that much difference won't be there.

Hence it is Robust to outliers.



② Disadvantage of MAE

① Convergence usually takes more time.

② Optimization is a complex task

* In MSE the graph comes as parabola since it gives quadratic eqn.



* But In Case of MAE, we have linear eqn. Hence the graph comes in straight line.

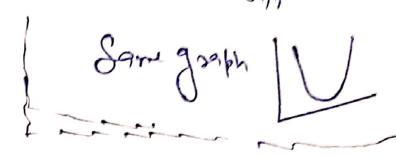
Here we find the derivative with the help of Subgradient Concept.
Subgradient basically says when we can't find the derivative at 0, divide into regions to find the derivative.

We can't find the derivative at 0, that you have to divide into regions by region and try to find out the derivative to come near 0.

Divide into parts by parts and find out in Subgradient and then take the derivative.

iii) Time Consuming:

Assignment

- ① Huber Loss → Combination of MAE & MSE. } Study yourself
some on MSE
- ② RMSE → \sqrt{MSE} → } Read about Advr. & Disadv.
- NOTE } W.r.t. Unit
- * If there is outlier we use MAE. } Outliers & Same MSE
- If there is no outlier we use MSE. } Differentiable
- ↑ diff.
- Same graph 
- # Performance Metrics : →

Q How do we check if the model is good or not in case of Linear Regression.

Loss function = MAE, MSE, RMSE, Huber Loss

→ for that basically we need to check the performance matrix.

- i) R Squared ii) Adjusted R Squared { It helps to check the performance of the model.

i) R Squared :

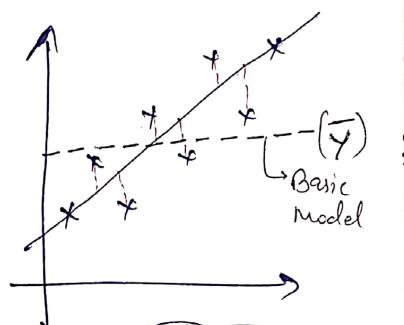
$$R^2 = 1 - \frac{SS_{Res}}{SS_{Total}}$$

where SS_{Res} = Sum of Square Residuals.

SS_{Total} = Sum of Square Avg.

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

where \bar{y} = Avg. of y .



If model is fitted well, numerator will be less since error will be less and denominator will be high.

$$= 1 - \frac{\text{Small no.}}{\text{Biggen no.}}$$

→ Small no.

$$= 1 - \text{Small no. i.e. } \leq 1$$

Hence R-Squared value will be always less than 1.

iff, R^2 Squared = 0.88 Model is

\rightarrow 85% Accurate.

iff, R^2 Squared = 0.75 Model is

\rightarrow 75% Accurate.

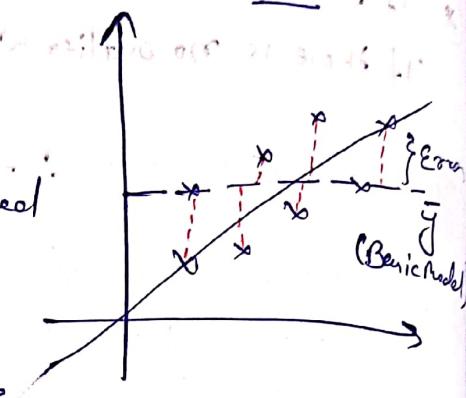
Q. Why denominator is high?

W.R.T \bar{y} Error will be more

and since it is not best fitted

Scatterline. Hence Squaring that

will make the error bigger.



④ R Squared measures the performance of the Model

that you have actually created.

Q. Can R^2 Squared value be -ve?

If at all R^2 Squared value is -ve it means
Model is very bad. Even bad than the \bar{y} .

Bez.

$$R^2 \text{ Squared} = 1 - \frac{\text{Big no.}}{\text{Small no.}}$$

Big no.
Small no.

If num. is big
means error is
big means very
very bad model

$1 - \frac{\text{Big no.}}{\text{Small no.}}$

Even bad than \bar{y}

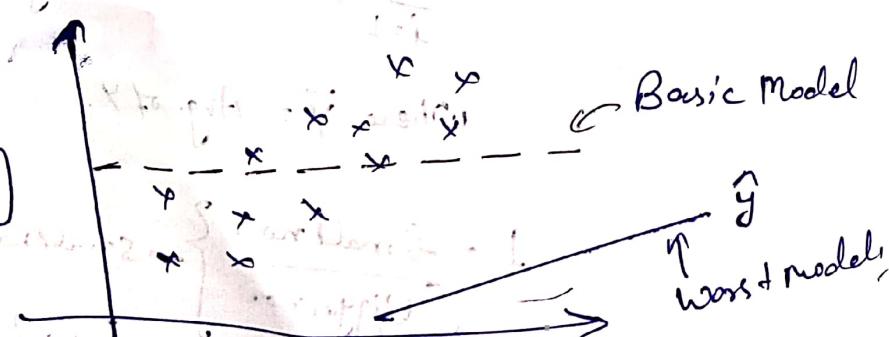
$$R^2 \text{ Squared} = -\text{ve} \text{ i.e. } < 0$$

Basic Model
e.g.

Eg:-

$$R^2 = -\text{ve.}$$

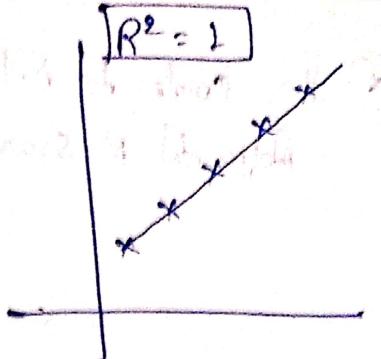
In This Case



④ If all the data points are coming over the best fit line, in that case $R^2 = 1$

Bcz Error = 0

$$\text{Hence } R^2 = 1 - 0 = \boxed{1}$$



ii) Adjusted R Squared

Size of house	City Location	No. of bedrooms	Gender	Price
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-

No Correlation

⑤ If we keep adding feature if it is correlated with our O/P value i.e Price then R-Squared value ↑↑ Increases since the Model Can predict better Price.

But if we add a feature 'Gender' in the above example which is not correlated with price in that case R-Squared Price still increases. To remove this problem we use Adjusted - R Squared

R^2	Accuracy.
	65%
Size of house	
Size of house, City Loc	75%
Size of house, City Loc, No. of bedrooms	88%
Size of house, City Loc, No. of bedrooms, Gender	90%

Whenever the features are highly correlated the increase in accuracy happens very quickly.

→ no correlation still slight increase in Accuracy.
→ But this shouldn't happen since Gender & Price is not correlated

* In order to solve this problem we use Adjusted R. Square.

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

Where, N = no. of data points

P = No. of independent factors

Feature	R^2	Adjusted R^2	It will be always be less than R^2 .
Size	65%	63%	$P=1$
Size, City	75%	73%	$P=2$
Size, City, No. of bedro	88%	-	-
Size, bed room, Gender	90%	88%	$P=3$

Increase in Number of feature leads to decrease in accuracy when feature is added.

* Adjusted R^2 will always be less than R^2 as we can see from above example. In this case, the formula

* In case of Adjusted- R^2 if the feature is not

Correlated i.e. Gender \rightarrow Price
feature is not correlated

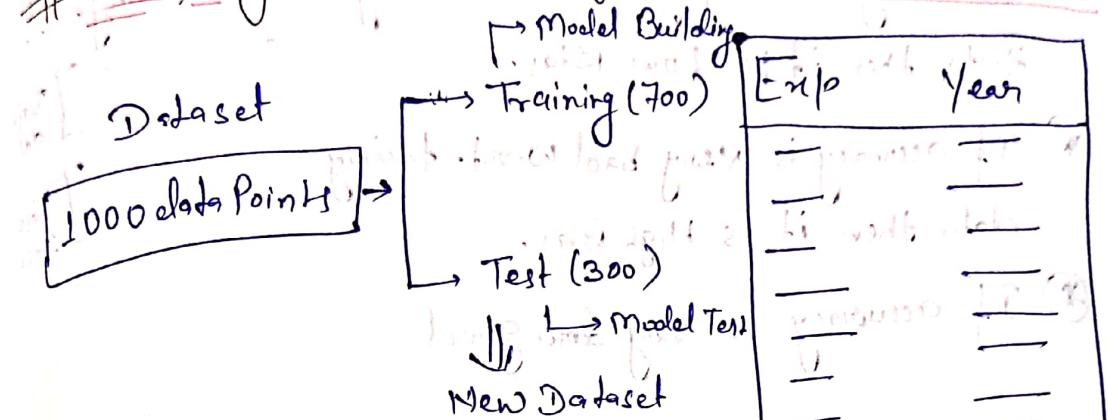
In that case Accuracy decreases.

* Means Adjusted- R^2 will not make any impact on Accuracy when unnecessary feature is added.

* Adjusted R^2 evaluates the accuracy based on only important features. Hence it is better than R^2 .

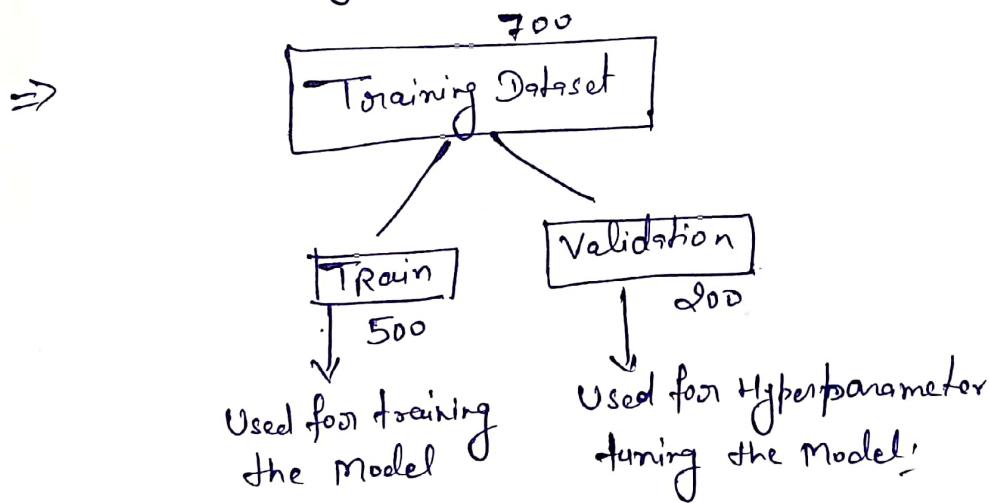
- ④ In Case of Performance Matrix we apply both R^2 and Adjusted R^2 .

Overfitting and Underfitting (Bias and Variance)



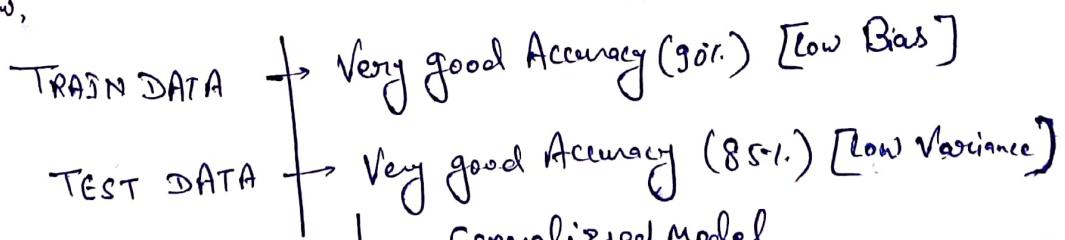
① Out of 1000 dataset suppose we are using 700 dataset for training and 300 dataset for testing.

② Training dataset has no idea about that 300 testing dataset. It is completely new for them



⇒ Suppose we have created one Model

Now,



③ The main of Model is to have Very good accuracy w.r.t TRAIN DATA and TEST Data and i.e. called Generalized Model.

③ Training data Accuracy is given by Bias and Test data Accuracy is given by Variance

- ④ If accuracy is very good w.r.t training data then it is Low Bias.
- ⑤ If accuracy is very bad w.r.t. training data then it is High Bias.
- ⑥ If accuracy is very bad w.r.t test data then it is Low Variance.
- ⑦ If accuracy is very bad w.r.t test data then it is High Variance.
- # Overfitting →

TRAIN DATA → Very good Accuracy [90%] [Low Bias]

TEST DATA → Bad Accuracy [50%] [High Variance]

↓ This condition is called

Overfitting → Solve using Hyperparameter Tuning

- ⑧ Overfitting is a situation wherein Model gives very good accuracy w.r.t Train data but gives bad accuracy w.r.t Test data.

Eg:- Ratta Man, Cracking Hugs without understanding

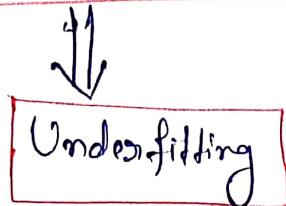
Q How do we solve this Problem?

By Performing Hyperparameter Tuning

We increase the data size after test

Underfitting :-

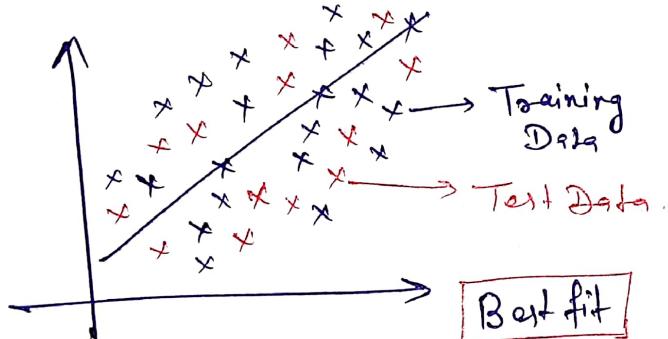
TRAIN DATA → Model Accuracy is Low [High Bias]
TEST DATA → Model Accuracy is Low/High [Low or High Variance]



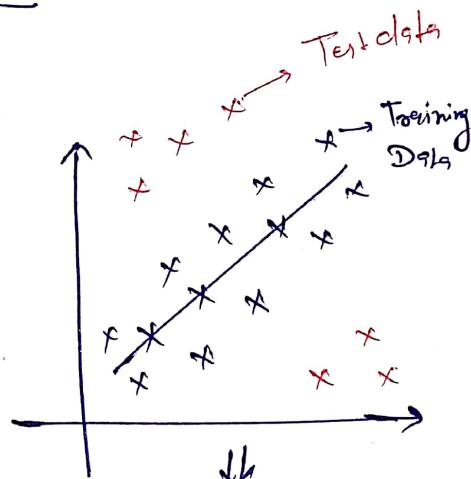
- ① When Model accuracy is very low w.r.t training data and if accuracy is high or low w.r.t test data that situation is called Underfitting.

Eg:- Students who doesn't study and In exam they just do tukka mas, So result could be either good or bad i.e unpredictable.

In Case of generalized Model

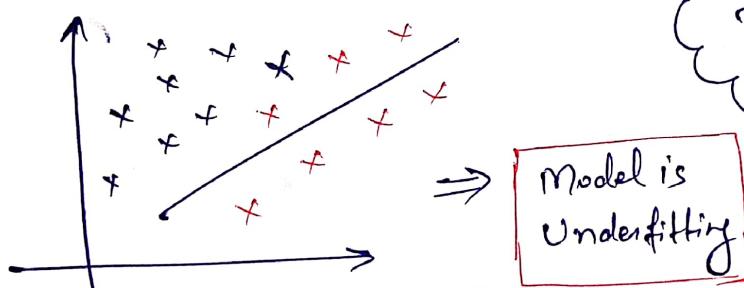


Model is underfitting,
Generalized model



Model is Overfitting

Trained well with
training data but
not with the test data



Model is
Underfitting

Test data may perform
good in this case