

ML - Day 2

Linear Regression Algorithm

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

Convergence Algorithm

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

↓ Predicted Pts. ↓ Truth Points

Cost Function Vs Loss Function → Every observation

The Loss function is to capture the difference b/w the actual and predicted values for a single record.

Cost function aggregates the difference b/w actual and predicted points for the entire training dataset.

$$\text{Loss function} = (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\Rightarrow (\hat{y}^{(i)} - y^{(i)})^2$$

↓ Predicted ↓ Actual
 pts. pts.

→ Calculating the slope

$$\text{ste } \theta_j : \theta_j - \alpha \boxed{\frac{d J(\theta_j)}{d \theta_j}}$$

derivative ↗

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

for $j = 0$

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

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

for $j = 1$

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

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

Repeat until Convergence

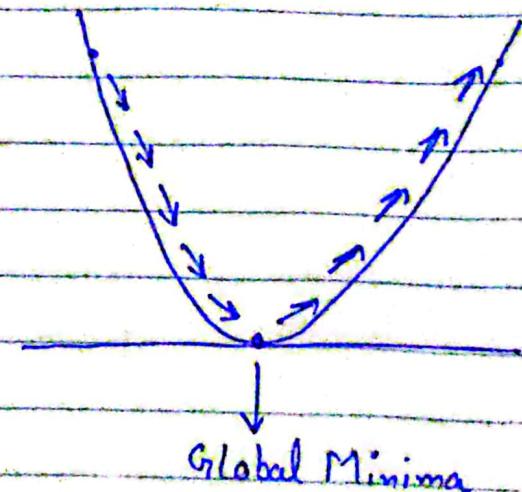
{

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

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

}

α = Speed of Convergence



Cost Functions

- ① Mean Squared Error (MSE)
- ② Mean Absolute Error (MAE)
- ③ Root Mean Squared Error (RMSE)

MSE (Mean Squared Error)

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

Quadratic E.qn

$$(a-b)^2 \downarrow \\ a^2 - 2ab + b^2 \downarrow \\ ax^2 + bx + c = 0$$

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

\hat{y} → predicted values.

When we

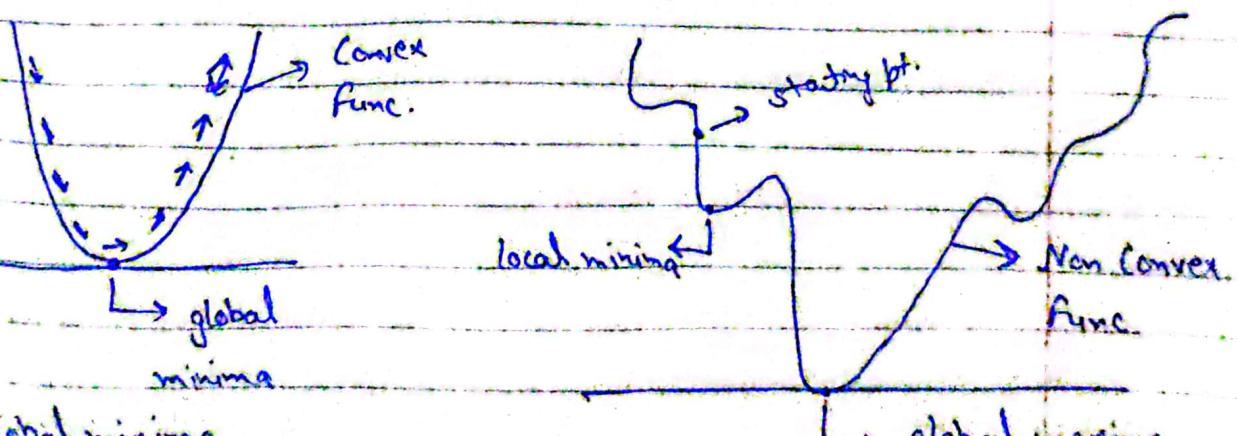
→ Advantages of MSE

① This equation is differentiable.

② The equation has only one global Minima.



When we plot quadratic egn we get a parabola



At global minima



slope = 0

In this case at
local minimum

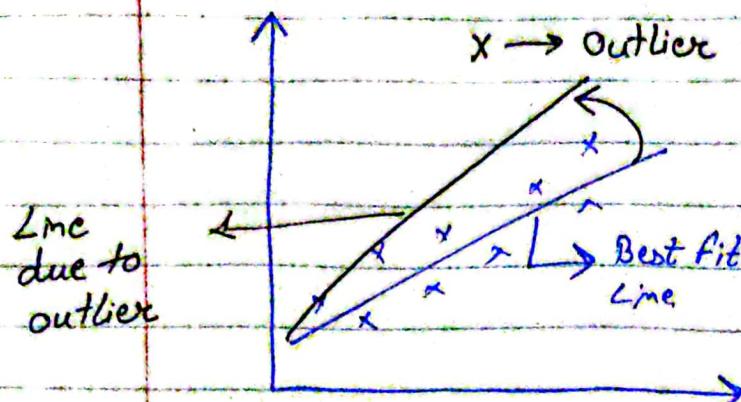


slope = 0, so convergence
will not go upto global minima

we won't be able to reduce
Cost function.

→ Disadvantages of MSE

- ① This is not robust to outliers.



Because of this outlier \rightarrow the cost function will increase by a huge amount

↓
Best fit line will get updated by a big margin

In order to solve this problem we need to remove the outliers.

- ② Penalizing The error \rightarrow Changing the Unit

Let say we have a data

Independent Exp. Salary (Lakhs ins)

$$(y - \hat{y})^2$$

$$(\text{salary} - \text{pred.salary})^2$$

(Error)²

↓ As the units for
penalized salary is Lakhs

↓ increasing the
order

↓ so when we
square the difference
we are also squaring Lakhs

Unit changes

$$(\text{Lakhs})^2$$

② Mean Absolute Error (MAE)

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

→ Advantages of MAE

① Robust to Outliers

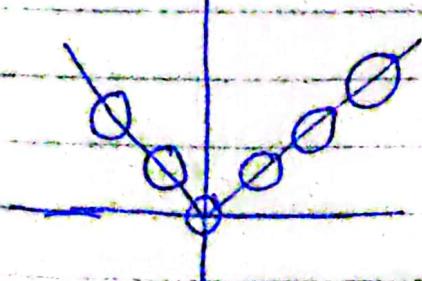
② It will also be in the same unit.

→ Disadvantages of MAE

① Convergence usually takes more time. Optimization is a complex task.

For MAE → we cannot find derivative at 0.

↓ To find derivatives we use Sub-Gradient concept.



③ Time consuming.

③ Huber Loss :-

→ Huber Loss is a combination of MSE and MAE.

→ We use Huber loss anytime we need a balance b/w giving outliers some weight, but not too much.

↓

for cases, where outliers are very important or there are not outliers we use MSE.

for cases, where we don't care at all about the outliers, we use MAE.

$$\text{Huber loss} = \begin{cases} \frac{1}{2} (y - \hat{y})^2, & |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2, & |y - \hat{y}| > \delta \end{cases}$$

Hyperparameter (δ), δ is used to switch two error functions.

→ Advantages of Huber Loss

① Outliers are handled properly.

② Local minima situation is handled here.

→ Disadvantages of Huber Loss

① It is complex.

② In order to maximize model accuracy, the hyperparameter δ will also need to be optimized which increases training requirements.

④ Root Mean Squared Error (RMSE)

$$\boxed{RMSE = \sqrt{MSE}}$$

→ Advantages of RMSE

It will also be in the same unit.

→ Disadvantages of RMSE

Not robust to outliers.

Performance Metrics



Helps to check the performance of the model.

① R-squared :-

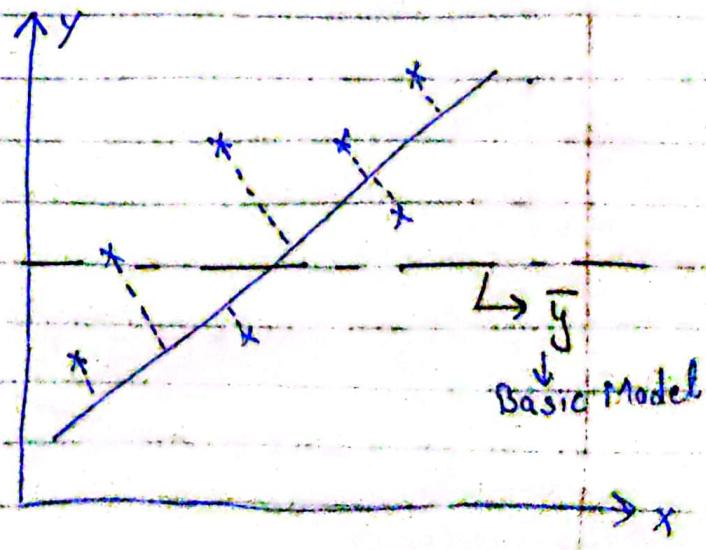
$$R\text{-squared} = 1 - \frac{SS_{Res}}{SS_{Total}}$$

SS_{Res} = Sum of Square Residuals
 SS_{Total} = Sum of Square Average

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

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

$$\bar{y} = \text{Average of } y$$



→ If model is fitted well $\rightarrow SS_{Res} \rightarrow$ Low & $SS_{Tot} \rightarrow$ High

$$R^{\text{Squared}} = 1 - \sum_{i=1}^n (y_i - \hat{y}_i)^2 \rightarrow \text{Low Value}$$

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

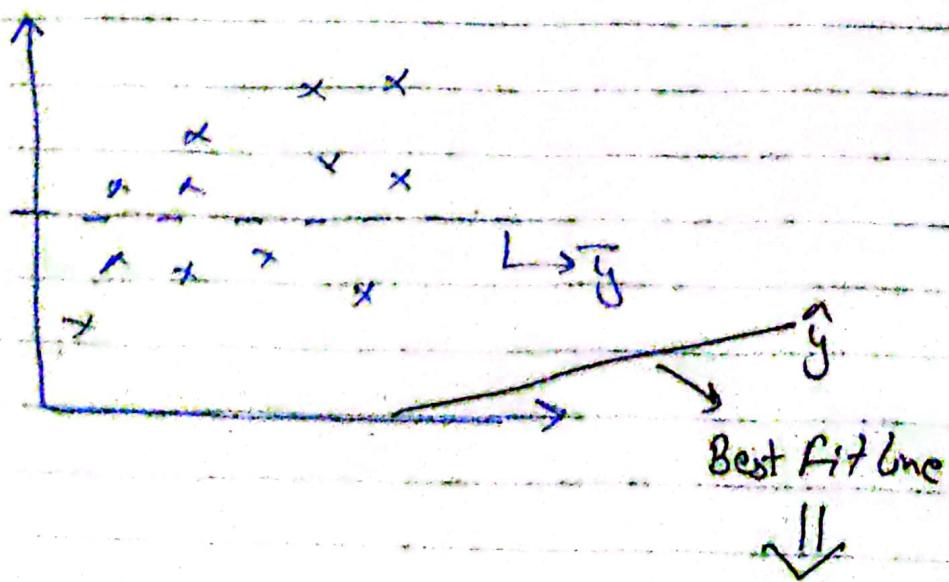
$$= 1 - \frac{\text{Small Num}}{\text{Bigger Num}} \rightarrow \text{Small Number}$$

$$\Rightarrow R^{\text{Squared}} < 1$$

Let $R^{\text{Squared}} = 0.85 \rightarrow 85\%$ accurate

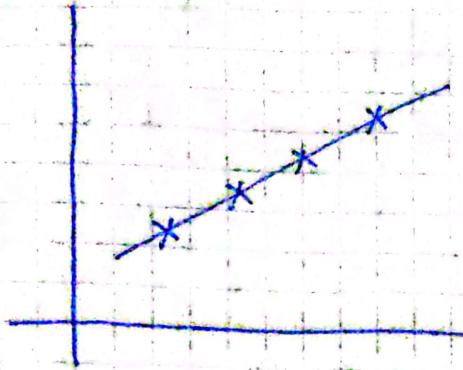
R-Squared measures the performance of model

→ If R-Squared value is -ve, then the model is
* very bad.



R-Squared is -ve

$$\rightarrow R^2 = 1$$



② Adjusted R-squared

Let say we have a dataset

Size of House	City Location	No. of Bedrooms	Gender	Price

Here, size of house, city location and num of bedrooms are highly correlated with the Price.

Gender is not highly correlated with Price

Now

① 2 features only Size & Price

$$R^2 = 65\%$$

② Add a feature to dataset, say city location

$$R^2 = 75\%$$

③ Add another feature, no. of bedrooms

$$R^2 = 88\%$$

④ Add gender column

$$R^2 = 90\%$$

} when features are Highly

} correlated, the accuracy will increase rapidly

} Accuracy will be very less increased.

Here the accuracy should not have increased. To solve this problem we use adjusted R^2 .

Because there is no direct correlation

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

N = No. of Data Points.

P = No. of Independent Features.

when P

when,

$$P = 1, R^2 = 65\% \rightarrow \text{Adj } R^2 = 63\%$$

$$P = 2, R^2 = 75\% \rightarrow \text{Adj } R^2 = 73\%$$

$$P = 3, R^2 = 88\% \rightarrow \text{Adj } R^2 = 87\%$$

$$P = 4, R^2 = 90\% \rightarrow \text{Adj } R^2 = 83\%$$



when we are adding Gender to our dataset, which is not correlated to Price column

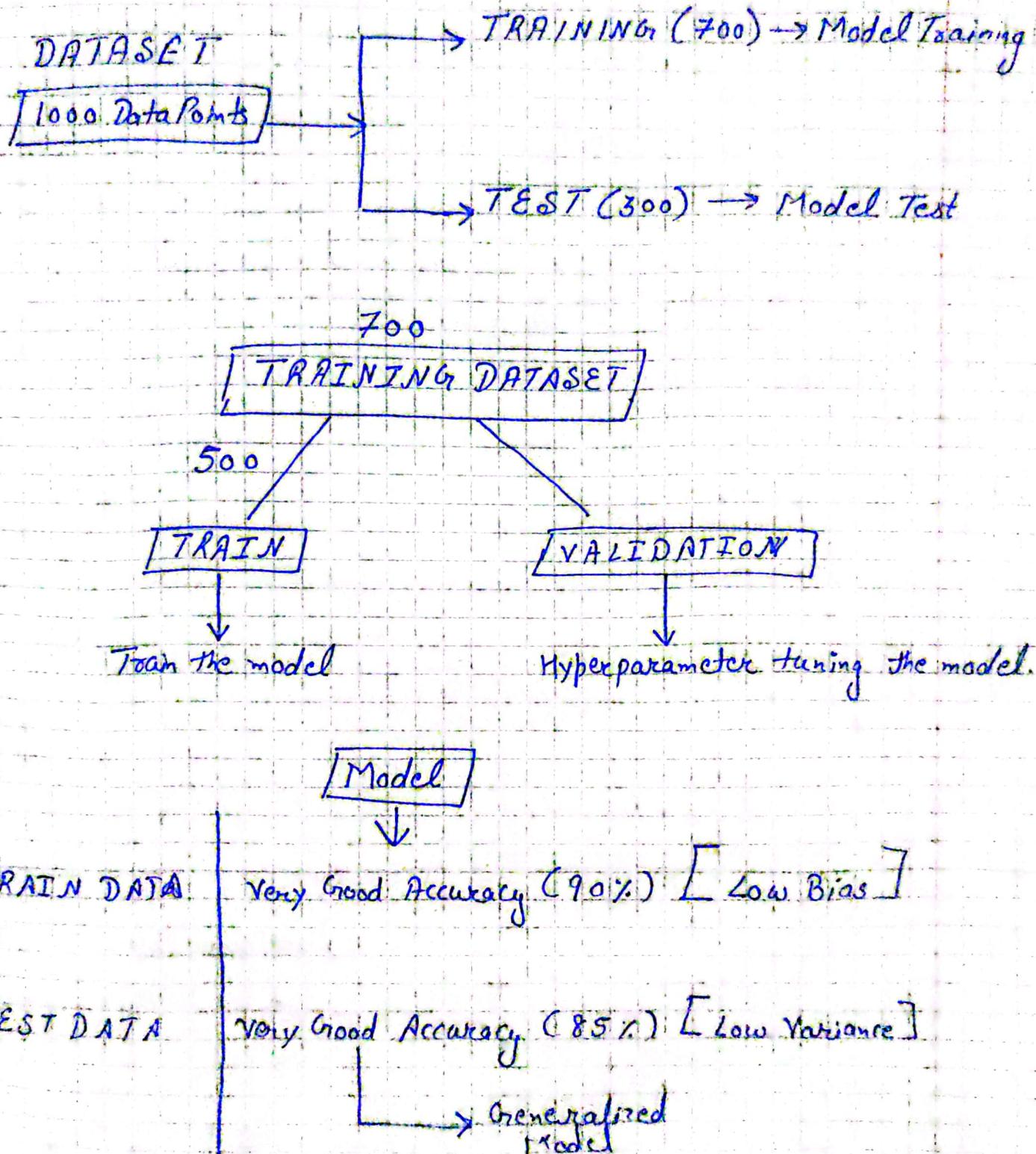


R^2 is increasing but $\text{Adj } R^2$ is decreasing

→ Adjusted R^2 will measure the performance of model based on very important features.

Overfitting & Underfitting

(Bias and Variance)



* Aim :- Our model should be able to get good training and good test accuracy.

\rightarrow TRAIN	Very Good Accuracy (90%) [Low Bias]
TEST	Bad Accuracy (50%) [High Variance]

↓

Overfitting

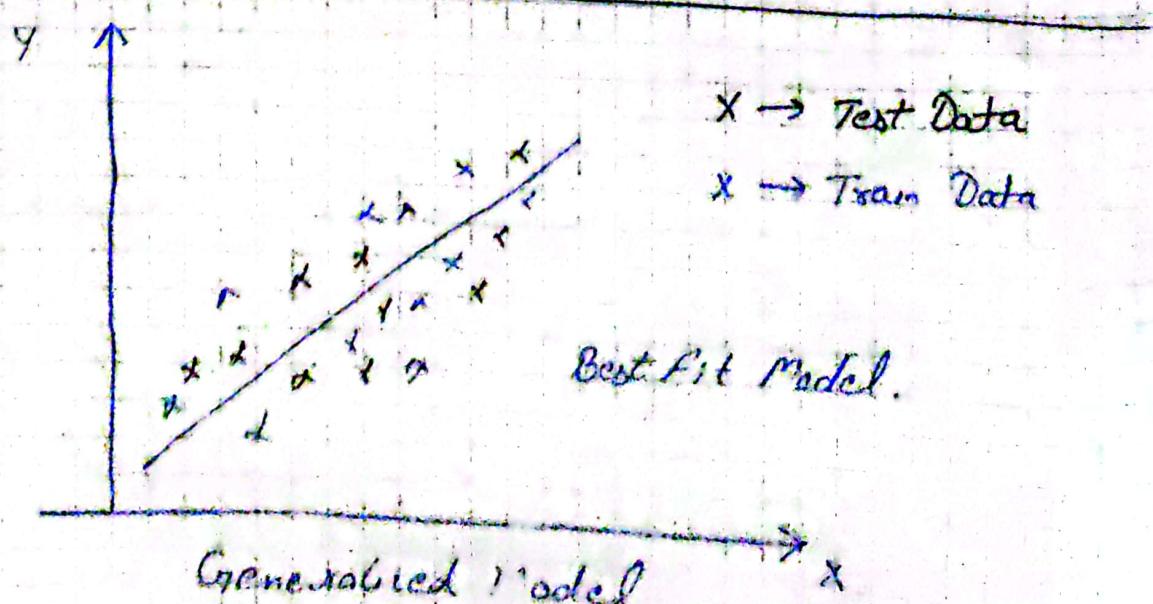
- \rightarrow To solve overfitting problem, we need to perform Hyperparameter Tuning
 or
 we can increase the data.

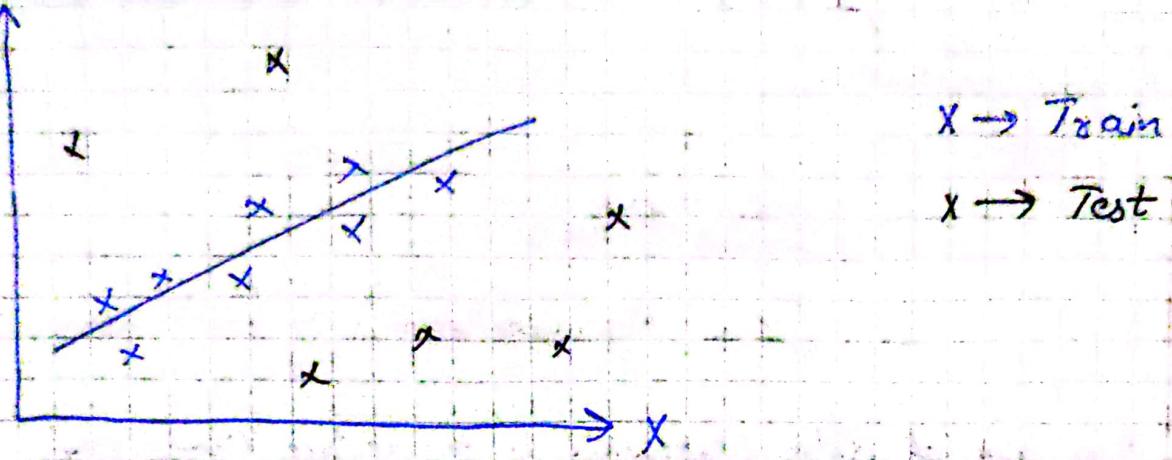
Now,

\rightarrow TRAIN	Model accuracy is low [High Bias]
TEST	Accuracy is Low / High [Low or High Variance]

↓

Underfitting





Overfitting

