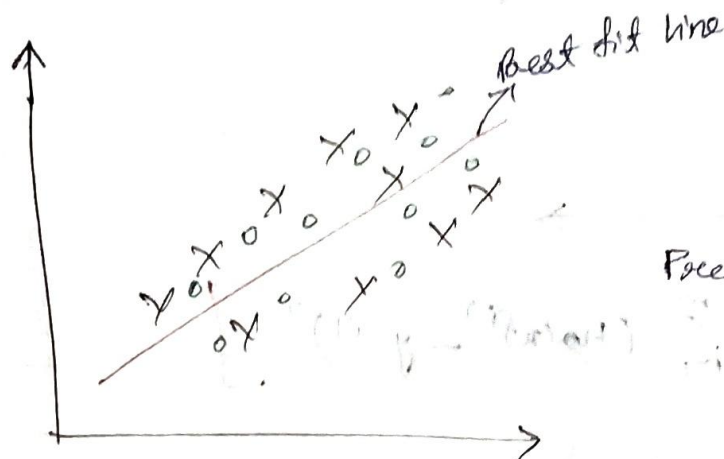


Underfitting

Train Data \rightarrow low accuracy \rightarrow [HIGH BIAS]

Test Data \rightarrow High Accuracy / low accuracy \rightarrow [HIGH / LOW VARIANCE]

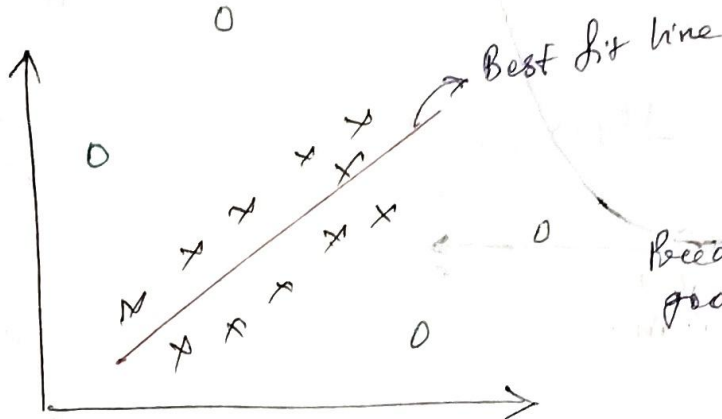
Generalize Model



x \rightarrow Training Data
o \rightarrow Test Data

Prediction will be good.

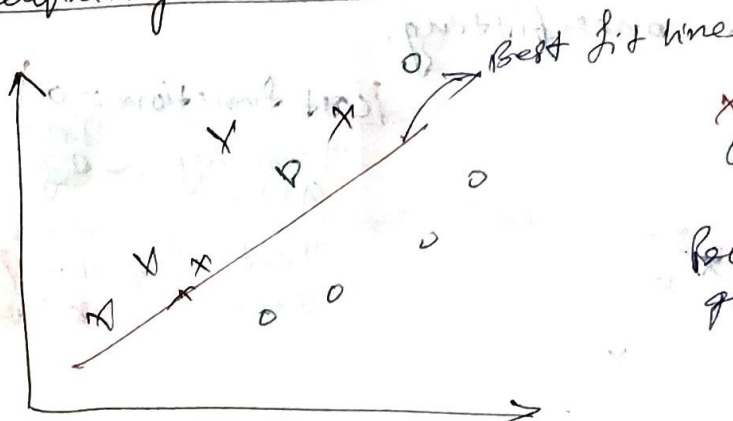
Overfitting Model



x \rightarrow Training Data
o \rightarrow Test Data

Prediction will not be good.

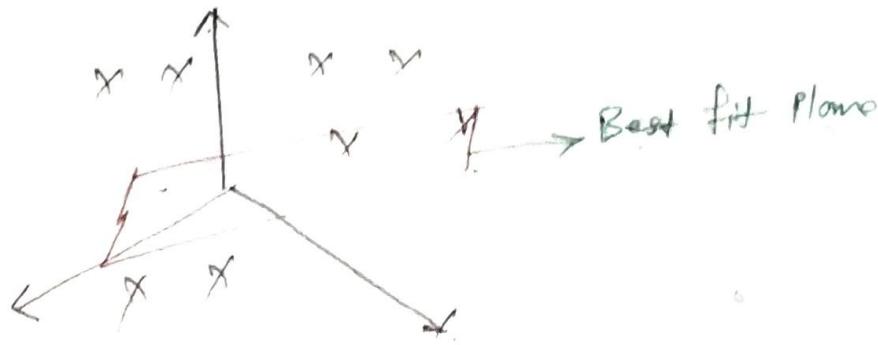
Underfitting Model



x \rightarrow Training Data
o \rightarrow Test Data

Prediction will not be good.

3-D



More than 3-D \rightarrow Hyperplane.

09-01-2022

M2 - Day - 2

Convergence Algorithm

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \quad \rightarrow \text{cost function}$$

Difference b/w cost function and loss function

cost function

- we find error for all the points and do average

loss function

- we find error for observed points and if we find error for all the points, then it is loss function.

$\text{loss function} = (h(x)^{(i)} - y^{(i)})^2$	
\downarrow	
Predicted value	Actual value

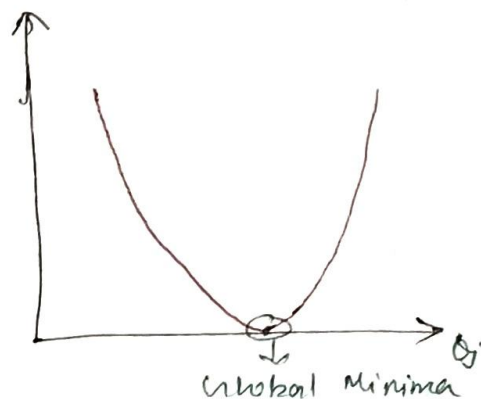
Repeat until convergence,

{

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

}

$J(\theta_0, \theta_1)$



→ To achieve global minima.

Derivative w.r.t θ_0 , $\boxed{J=0}$

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

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

Just for derivative purpose.

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

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

$$\begin{cases} \frac{\partial \theta_0}{\partial \theta_0} = 1 \\ \frac{\partial \theta_1 x}{\partial \theta_0} = 0 \\ \frac{\partial y^{(i)}}{\partial \theta_0} = 0 \end{cases}$$

Derivative w.r.t θ_1 , $\boxed{J=1}$

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

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

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

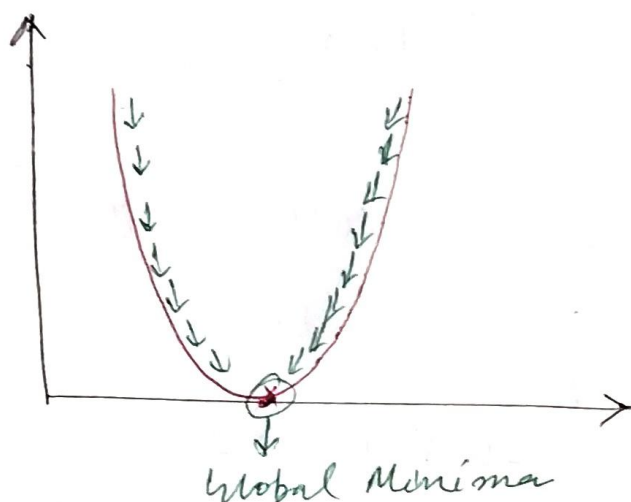
$$= \frac{1}{m} \sum_{i=1}^m \{(\theta_0 + \theta_1 x)^{(i)} - y^{(i)}\} \times x$$

Repeat until convergence,
 $\{$

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

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

$\}$



Learning Rate: Speed of convergence.

Types of cost function

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

① MSE

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

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

↓
Predicted

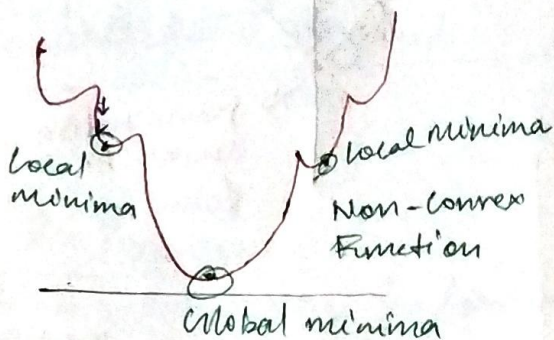
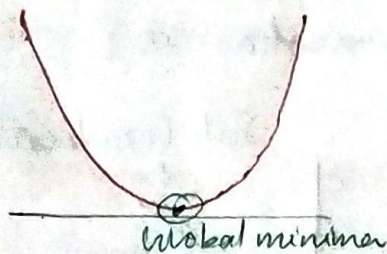
Advantages

- ① This equation is differentiable, we can calculate θ_0 and θ_1 .
- ② This equation also has only one global minima.

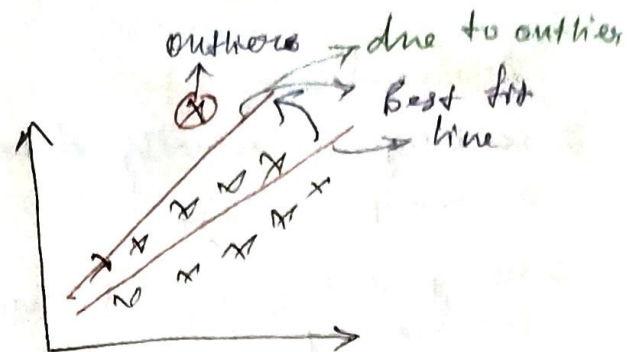
Disadvantages

- ① This equation is not robust to outliers.

convex Function



③ This equation has convex function.



→ Best fit line gets updated with huge margin in the presence of outliers.

Independent dependent
② $\text{Exp Salary (lakh INR)}$ $(y - \hat{y})^2 (\text{lakh})^2$

Time complexities
not increased

unit changing

we don't do scaling for dependent feature.

$(y - \hat{y})^2 \rightarrow \text{Squared} \rightarrow \text{Error} \rightarrow \text{penalized.}$
increased

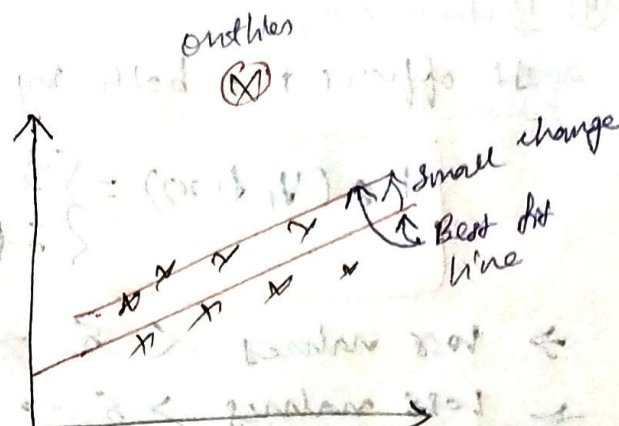
② MAE

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

$\hat{y} = 0.0 + 0.1x$
Predicted.

Advantages

① Robust to outliers.

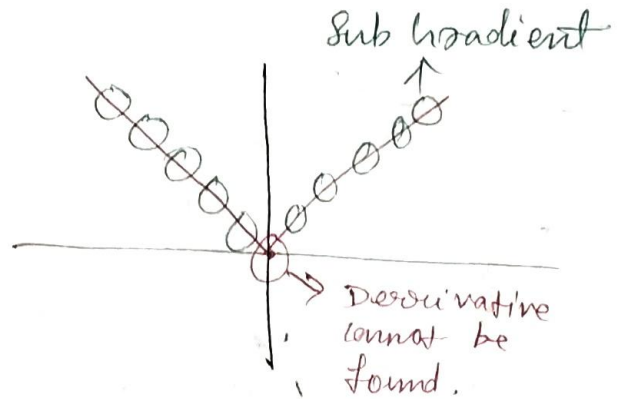


Disadvantages

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

- Using sub-gradient concept, we can find derivative.

② Time consuming.



Huber Loss

③ RMSE

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

Advantages

① It is differentiable.

② Unit will be same.

Disadvantages

① Not robust to outliers

④ Huber Loss

- It offers the both by balancing MSE and MAE together.

$$\boxed{L_{\delta}(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & , \text{for } |y - f(x)| \leq \delta, \\ \delta |y - f(x)| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}}$$

→ loss values $< \delta \rightarrow$ Use MSE

→ loss values $> \delta \rightarrow$ Use MAE

Note: Use the Huber loss any time you feel that you need a balance b/w giving outliers "some weights" but not too much.

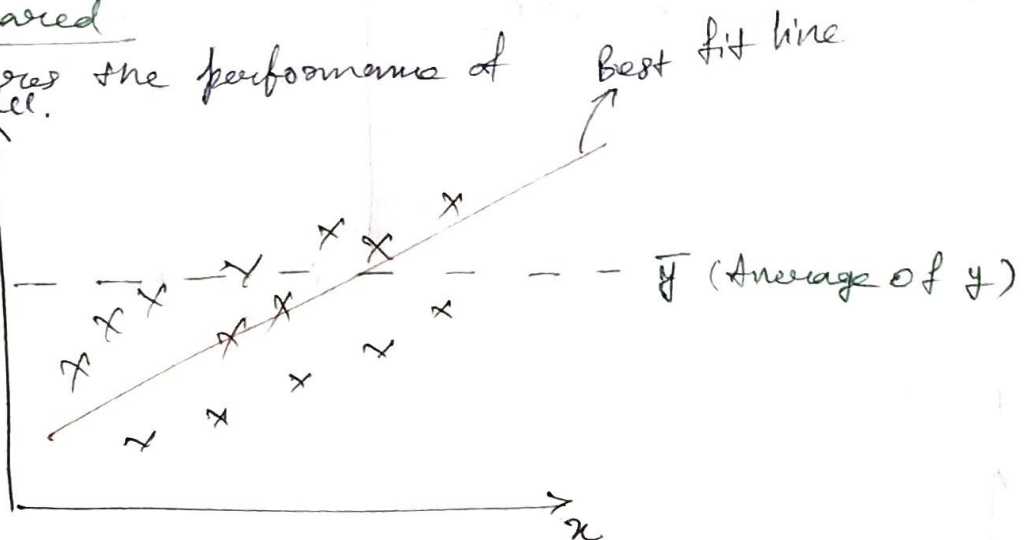
How can we check if a model is good or not?
- Using performance metrics

Performance Metrics

- ① R-Squared
- ② Adjusted-R Squared

① R-Squared

→ It measures the performance of the model.



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

SS_{Res} = Sum of Square Residuals

SS_{Total} = Sum of Square Average

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

→ Low value } for good model
→ high value }

\bar{y} = Average of y

$$= 1 - \left\{ \frac{\text{Small number}}{\text{Bigger number}} \right\} \text{ small number}$$

$$R\text{-Squared} \leq 1$$

If $R\text{-Squared} = 0.85 \rightarrow 85\%$ Accurate

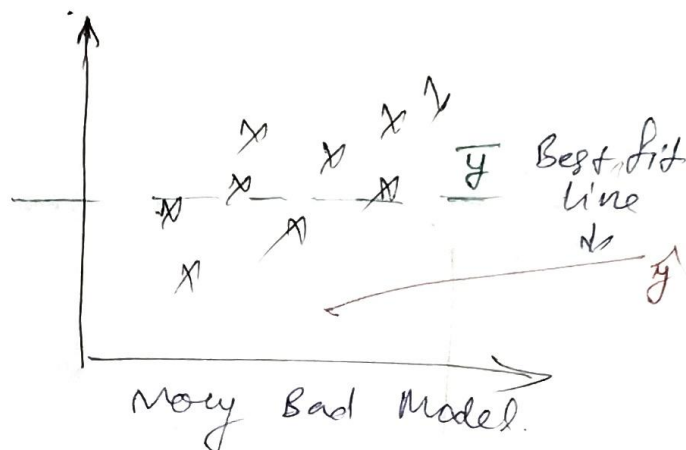
$R\text{-Squared} = 0.75 \rightarrow 75\%$ Accurate

Ques Can $R\text{-Squared}$ value be negative?

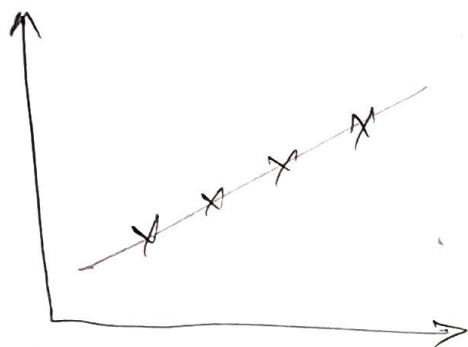
Ans If it is '-ve', then our model will be very bad.

Here $(y_i - \hat{y}) > (y_i - \bar{y})$

then, $R\text{-Squared}$ will be '-ve'.



If $R^2 = 1$



② Adjusted $R\text{-Squared}$

No direct correlation

Price

Price

Size of House City location ~~Price~~ No. of Bedroom Salary of the owner

\rightarrow Earlier only size of house was present as an independent feature to get the price and has some $R\text{-Squared}$ value. But with the inclusion of city location, $R\text{-Squared}$ value got increased. After that again with the inclusion of city location ~~and~~ no. of bedroom and hence, the value got increased. But as there is no direct correlation b/w breeder and price, the value shouldn't increase.

R^2	features inclusion	Adjusted R^2	Independent features
65%	Size of House	63%	$P=1$
75%	City location	73%	$P=2$
88%	No. of Bedrooms	85%	$P=3$
90%	Render.	85% (decreased)	$P=4$

→ No direct correlation w.r.t price.

very slight increase, but this shouldn't happen. To solve this problem, we use Adjusted R -Squared.

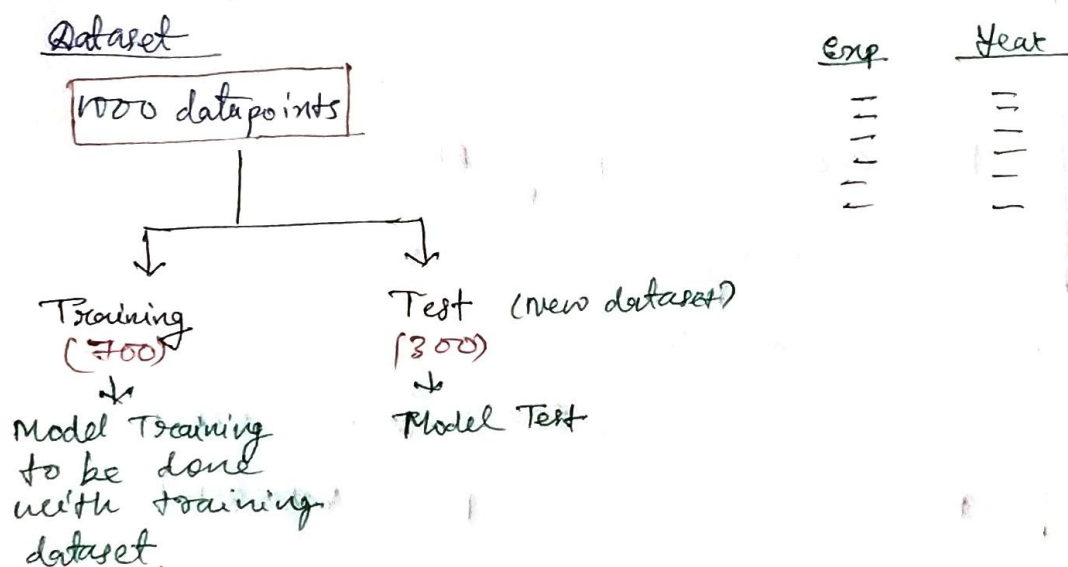
$$\text{Adjusted } R^2 = 1 - \frac{(1 - R^2)(N - 1)}{N - P - 1}$$

N = No. of data points

P = No. of independent features.

→ Adjusted R^2 is the better metrics to evaluate the model.

Overfitting and Underfitting (Bias and Variance)



Training Dataset (1000)

↓
TRAIN
(500)
↓
used for
training of
the model

↓
VALIDATION
(500)
↓
used for hyperparameter
tuning of the model

MODEL

Generalize Model

TRAIN DATA	→	Very good accuracy (90%)	} our aim is to get good train and test accuracy.
TEST DATA	→	Very good accuracy (85%)	

Train Data → very good accuracy (90%) → ^{Low} [BIAS]

Test Data → very good accuracy (85%) → ^{Low} [VARIANCE]

Train Data → Low accuracy → [High BIAS]

Test Data → Low accuracy → [High VARIANCE]

~~Overfitting~~ { Overfitting }
~~Very good~~

Train Data → Very good accuracy (90%) → [Low BIAS]

Test Data → Bad accuracy (50%) → [HIGH VARIANCE]

→ We can solve this by issue by performing
hyperparameter tuning or by increasing the no. of
dataset.