

Loop function up cost function

In case of cost function all the point are considered.

In case of loss function it is for every observation

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

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

Predicted Actual

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

→ cost function
→ MSE

→ convergence Algorithm

repeat until ~~Algorithm~~ convergence

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

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

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

$$\underline{j=0}$$

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

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

$$\underline{j=1}$$

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

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

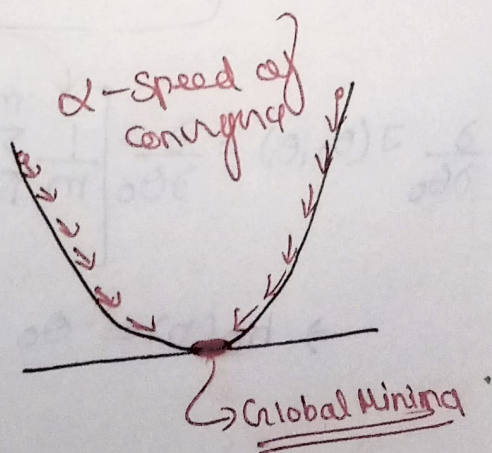
Repeat until convergence

{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m \left(h(\theta) - y^{(i)} \right)$$

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

}



Cost Function) ① MSE ② MAE ③ RMSE
(L1 loss)

① MSE {Mean Squared Error} (use when no outliers or to cast them)

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

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

→ Predicted value

→ Quadratic Eqⁿ Plot → Parabola

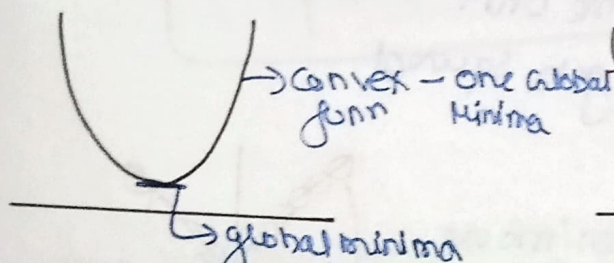
Advantage

① This Eqⁿ is Differentiable

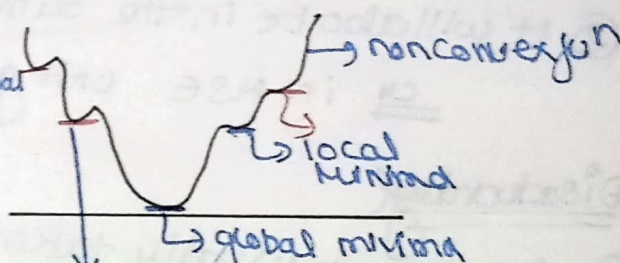
- Diff. needed to calc slope & update θ_0 & θ_1

② This Eqⁿ also has one global minima

always get parabola not having local minima.



• no local minima here slope = 0
∴ θ_0 & θ_1 not get updated



slope here is also '0'
∴ θ_0 & θ_1 will not get updated
get stuck at local minima & haven't reached global minima!

Disadvantage

① This is not robust to outliers

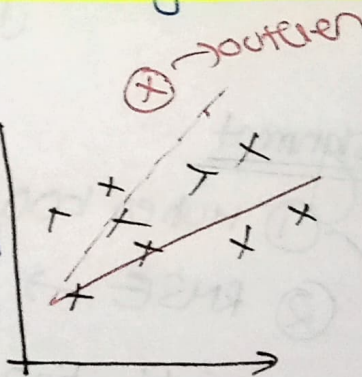
• coz of outliers → cost function

∴ Remove outliers update θ_0 & θ_1

② Normalizing the error (changing the unit)

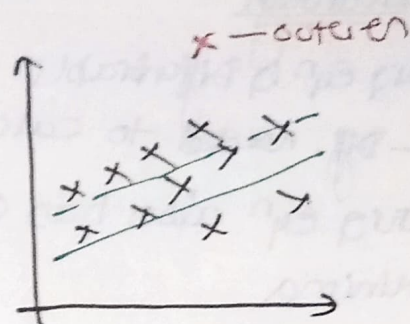
$$\text{Eg } (\text{real salary} - \text{predicted salary})^2 \text{ in Lakh}$$

$$= \text{Error}^2 \text{ (Lakh)}^2$$



② Mean Absolute Error (can be use with outliers)

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

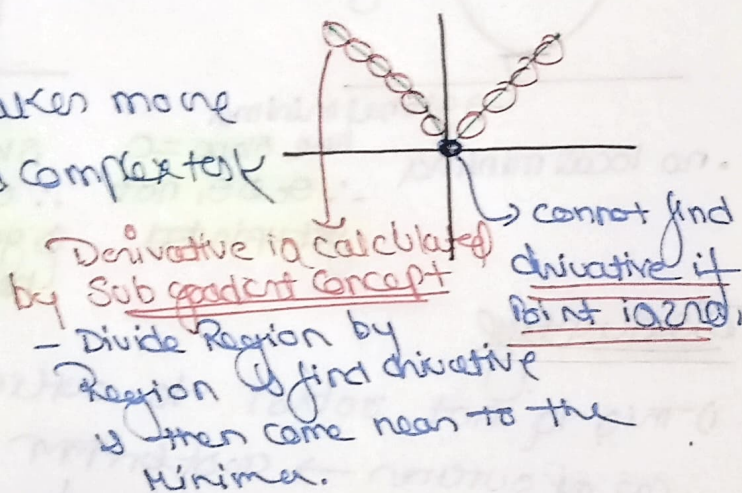


Advantage

- ① Robust to outliers
- ② It will also be in the same unit
ex in MSE unit gets squared

Disadvantage

- ① Convergence usually takes more time, optimization is a complex task
- ② Time consuming



Assignment

① Huber loss

② RMSE $\rightarrow \sqrt{MSE} \Rightarrow$ Average disadvantage

\rightarrow combination of MSE & MAE

- ③ Huber loss (Smooth Mean absolute Error)
- Combination of MSE & MAE. It takes good properties of both the loss function by being less sensitive to outliers & differentiable at minimum.
 - when error is less, MSE part of Huber is utilized & when error is large, the MAE part of Huber loss is used.
 - A hyperparameter is used to tell loss function to switch from MSE to MAE.

$$\text{loss} = \begin{cases} \frac{1}{2} * (n-y)^2 & \text{if } |n-y| \leq \delta \\ \delta * |n-y| - \frac{1}{2} * \delta^2, & \text{otherwise} \end{cases}$$

④ Root Mean Squared Error (RMSE)

$$L = \sqrt{\frac{1}{N} \left(\sum (\hat{y} - y)^2 \right)}$$

Other loss functions in ML:

1. Sum of Errors $L = \sum (\hat{y} - y)$
2. Sum of absolute Errors $L = \sum (|\hat{y} - y|)$
3. Sum of Squared Errors, $SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$
4. Mean Bias Error (MBE)

$$MBE = \frac{1}{n} \sum_{i=1}^n (p_i - o_i)$$

\hat{y} - predicted
 y - actual value

Performance Matrix

R Squared - Measure performance of Model that have been created.

$$R \text{ Squared} = 1 - \frac{SS_{\text{res}}}{SS_{\text{total}}}$$

SS_{res} = Sum of squared Residual

SS_{total} = Sum of squared avg.

$$= 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2 \rightarrow \text{low value}}{\sum_{i=1}^n (y_i - \bar{y})^2 \rightarrow \text{high value}} \left. \vphantom{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}} \right\} \text{if model fitted well}$$

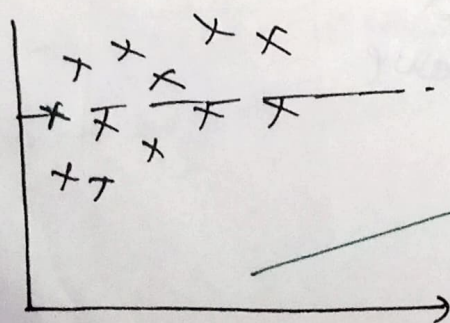
$$= 1 - \left\{ \frac{\text{small no.}}{\text{Bigger no.}} \right\} \rightarrow \text{small no.}$$

$R \leq 1$ - if model fitted well

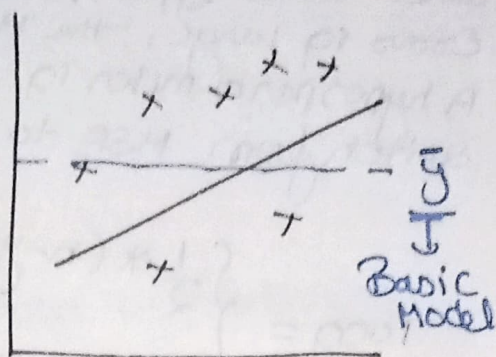
$\Rightarrow 0.85 \rightarrow 85\%$ accurate

$\Rightarrow 0.75 \rightarrow 75\%$ accurate

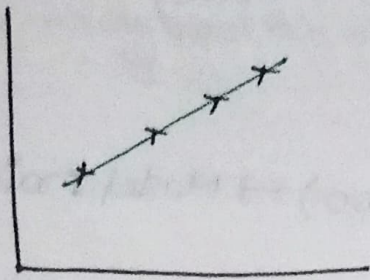
Q Can R squared value be -ve \rightarrow Model is very Bad



$$\underline{\underline{R^2 = -ve}}$$



* $R^2 = 1$



Adjusted R Square → will going to evaluate accuracy based on very ^{Dependent} Important features _{feature}

Size of House	City	No. of Rooms	Gender	Price
------------------	------	-----------------	--------	-------

- Coz of addition of more independent features the R^2 may inc which is Desired. But there may be such feature which have no correlation b/w feature & output but there may be inc in accuracy but inc in accuracy is very slight. In order to solve such problem adjusted R^2 is used.

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

N = no. of data points

P = no. of independent features

ex

R^2 65% $\underline{P=1}$

R^2 75% $P=2$

R^2 88% $!$

R^2 88 $\underline{P=4}$

Dependent
Added

→ less correlated
features

Adj $R^2 = 63\%$ } less than
Adj $R^2 = 73\%$ } R^2

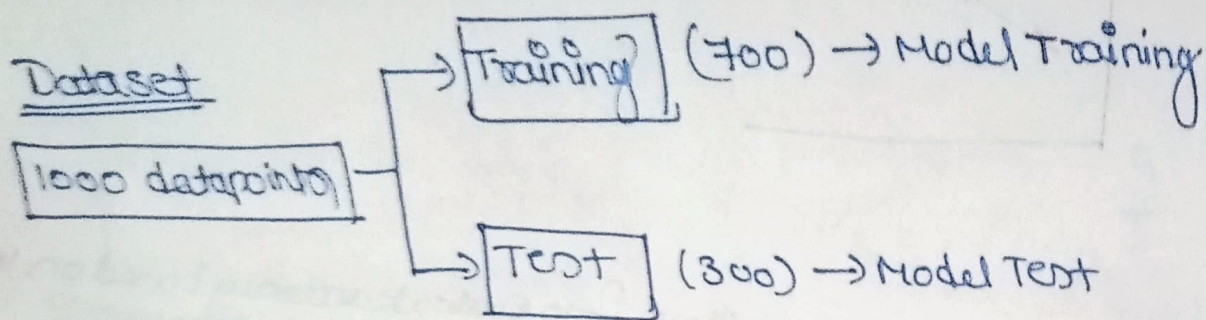
Adjusted $R^2 = 85\%$

↓
Dec adjusted R^2

Coz that feature is
not correlated.

- In inc unnecessary columns R^2 will inc but Adj R^2 will dec.

Overfitting & underfitting (Bias and Variance)



Training Dataset (700)

(500)

TRAIN

Train the Model

200

Validation

Hyperparameter tuning the Model

What's say

- For a Model

Train Data - very good Accuracy (90%) [Low Bias]

Test Data - Very good Accuracy (85%) [Low Variance]

Train Data Accuracy

_____ Bias

Test Data Accuracy

_____ Variance

*

TRAIN - very good Accuracy [90%] [Low Bias]

Test - Bad Accuracy [50%] [High variance]

⇓

overfitting

*

TRAIN - Model Accuracy is low

[High Bias]

TEST - Model Accuracy is low/high [low or high variance]

⇓

Underfitting

Model is underfitting

