



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

To achieve global minima:derivative w.r.t. $\theta_0, j=0$

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

for $j=0$ $\Rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x$

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

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

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

derivative w.r.t $\theta_1, j=1$

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

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



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

Repeat until convergence:

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

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

}

* learning rate:- speed of convergence

Types of cost function

① MSE :- Mean Square Error

② MAE :- Mean Absolute Error

③ RMSE :- Root mean Square Error

Mean Square Error:

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

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

↑
predicted



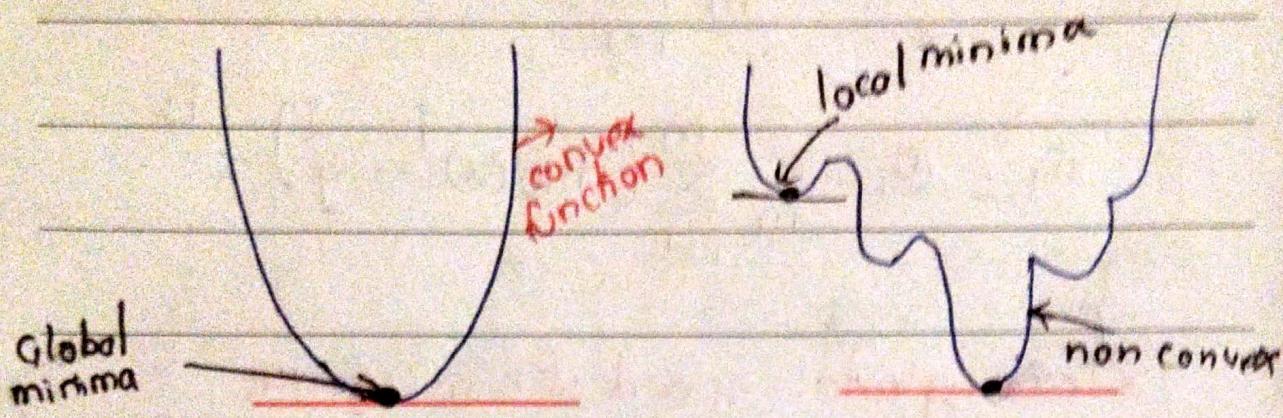
Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

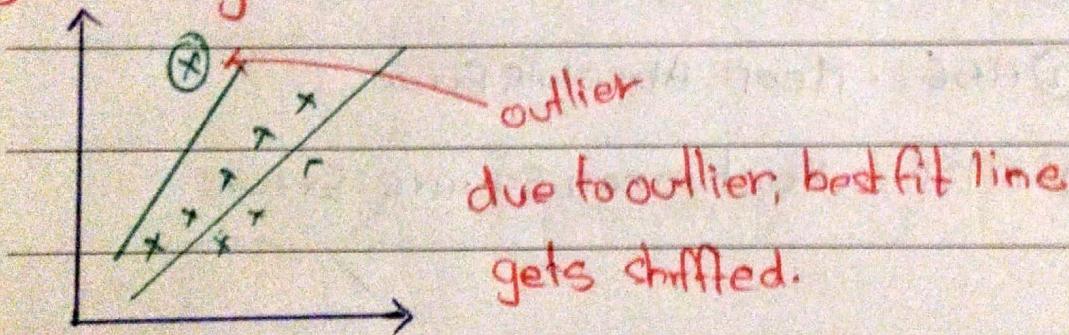
Advantage:- ① This equation is differentiable.

② It has only one global minima.



∴ main aim = it has convex function.

Disadvantage:- 1



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

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

Unit² → unit changing

time complexities increasing



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

We won't do scaling for depending features

$(y - \hat{y})^2 \Rightarrow$ got squared \Rightarrow error \Rightarrow penalized
value increase.

Mean Absolute Error

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

Advantage:

① Robust to outliers

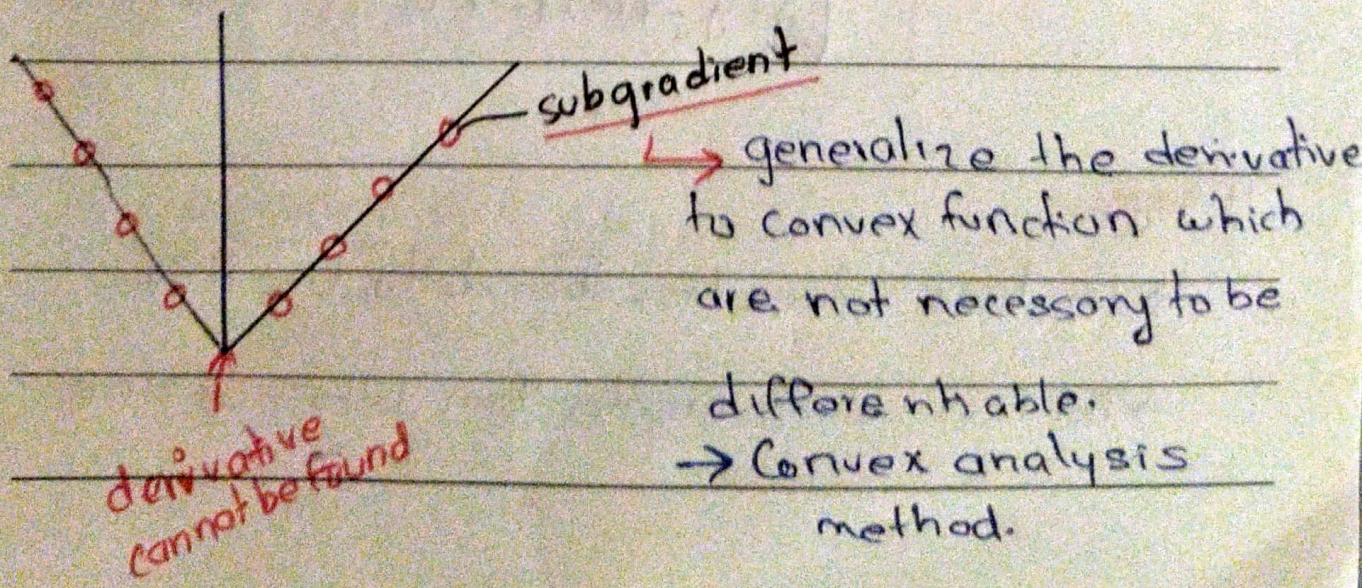
② It will be in the same unit

Disadvantage:

① Convergence usually takes more time.

Optimization is a complex task.

② time consuming





Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

Root Mean Square Error (RMSE)

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_i^n (y - \hat{y})^2}$$

Advantage

- It is differentiable.
- unit remains same

Disadvantage

- Not robust to the outliers

Huber loss function:

The huber loss offers the best of world by balancing the MSE and MAE together.

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

It says:

→ for loss values less than delta, use MSE

→ for loss values greater than delta, use MAE.



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

* MAE for larger loss

Using MAE for larger loss values mitigate

the weight that we put on outliers so that
we still get a well-rounded model.

MSE for smaller loss Why?

→ to maintain a quadratic function near
the centre.

Note:- Use the huber loss any time you feel
that you need a balance between giving
outliers some weight 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



Mo Tu We Th Fr Sa Su

Memo No. _____

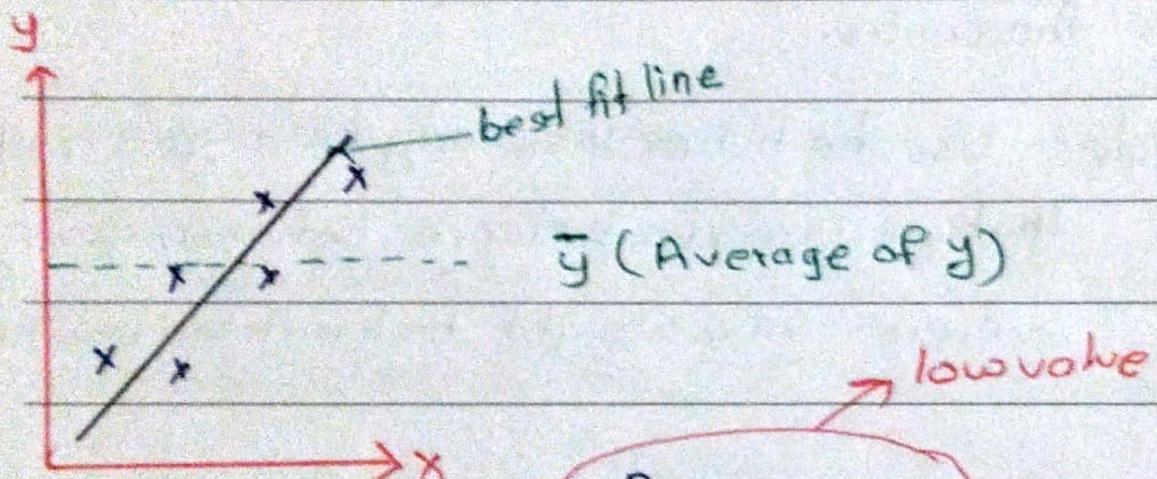
Date _____

3) R-Squared :- Measure the performance of the model.

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

SS_{RES} = sum of square residuals

SS_{Total} = sum of square average



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

high value

$$= 1 - \frac{\text{small value}}{\text{large value}} \rightarrow \text{small value}$$

$$\therefore R\text{-squared} \leq 1$$

Disadvantage } Even additional input
of R^2 : } variable shows no rels. the
variable, O/p variable,
Memo No. _____ Date _____ / R squared 1

Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

-- If R^2 of the model is 0.50, then approximately

half of the observed variation can be explained by the model's input.

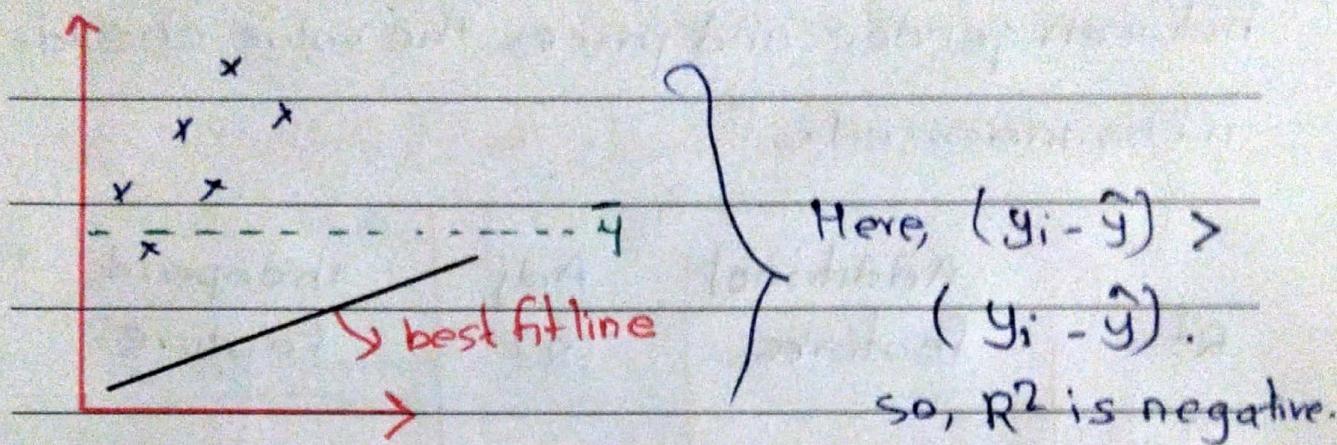


Fig: when $R^2 = -ve$.

→ model is not working good.

2) Adjusted R-Square

Eg: Model to predict Price of House.

size of house	city location	No of rooms	Gender	Price
---------------	---------------	-------------	--------	-------

No direct correlation

Before, when only size of house was present, as an independent feature it has some R-squared value. After, addition of features like city



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

location No. of bedrooms R^2 value got increased.
After addition of Gender as well, the R^2 value increased but there is no direct correlation between gender and price, the value shouldn't be increased.

R^2	Additional features	Adj R^2	Independent feature
65%	size of house	63%	$P=1$
75%	location	73%	$P=2$
88%	No. of bed	86%	$P=3$
90%	Gender	85%	$P=4$

↓
No direct relation with price

very slight increase, but this shouldn't

be happen. To solve this, we use Adjusted R^2 .



Memo No. _____

Mo	Tu	We	Th	Fri	Sa	Su
----	----	----	----	-----	----	----

Date / /

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

N = No. of data point

p = No. of independent Features

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