

09/10/22

Linear Regression Algorithm:

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

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

Convergence Algorithm: $\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} [J(\theta_j)]$

Annotations:
 α → step size
 $\frac{d}{d\theta_j} [J(\theta_j)]$ → derivative / slope
 $J(\theta_j)$ → cost function
 θ_j → here we are doing the summation

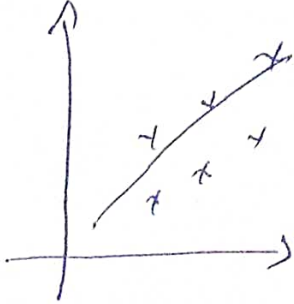
MSE

predicted point

truth point

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

[loss function vs cost function]
 ↓
 every observation.



loss function: $(h_{\theta}(x^{(i)}) - y^{(i)})^2 = (\underbrace{\hat{y}_0}_{\text{Predicted value}} - \underbrace{y^{(i)}}_{\text{Actual value}})^2$

Convergence : Repeat until convergence

$$\left\{ \begin{array}{l} \theta_j = \theta_j - \alpha \left[\frac{\partial}{\partial \theta_j} J(\theta_j) \right] \end{array} \right. \quad \text{learning curve}$$

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

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

$$\frac{\partial}{\partial x} (x+1)^2 = 2x(x+1) \times (1+0)$$

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

if $J=0$: $h_{\theta}(x) = \theta_0 + \theta_1 x$

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

$$= \frac{\partial}{\partial x} \frac{1}{2m} [(1+x) - y^{(i)}]^2 \quad \rightarrow \text{let us consider } (1+x) = \theta_0 + \theta_1 x$$

$$= \frac{2}{2m} [(1+x) - y^{(i)}] \times 1$$

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

$$\frac{\partial}{\partial \theta_0} (\theta_0 + \theta_1 x) = 1$$

if $J=1$:

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

$$= \frac{2}{2m} (\theta_0 + \theta_1 x - y) \times [x]$$

learning curve \rightarrow speed of convergence

Repeat until convergence

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \end{array} \right.$$

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

}

Cost function :-

1. MSE 2. MAE 3. RMSE 4. Huber loss.

Root Mean Squared Error

MSE = Mean Squared Error

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

$\hat{y} = c_0 + c_1 x$
 \hat{y} = predicted value

quadratic equation $\Rightarrow a^2 + b^2 = c^2$

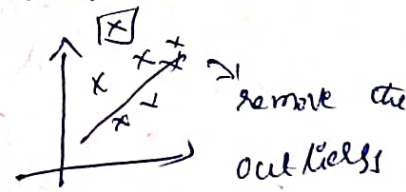
Advantage

- 1. This equation is differentiable
- This also has one global minima.



Disadvantage

- This equation is not robust to outliers



- Penalizing the error \Rightarrow increasing

- There is no local minima.

Exp	Salary (lacs) INR
-	-
-	-
-	-
-	-

$(y - \hat{y})^2$ squared \rightarrow Error \rightarrow Penalized

$(\text{Salary} - \text{Predicted Salary})^2 \times (\text{lacs})^2$

MSE $\Rightarrow \downarrow \downarrow \downarrow$

2. Mean Absolute Error (MAE)

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



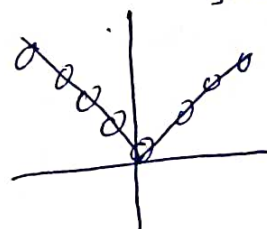
Advantage

1. Robust to outliers
 2. It will also be in the same unit.
- here we are not squaring the value

Disadvantage

- Convergence usually takes more time - optimization is a complex task

Sub Gradient Concept



- Time consuming.
- There may be a local minima.

Huber Loss

Combination of MSE & MAE

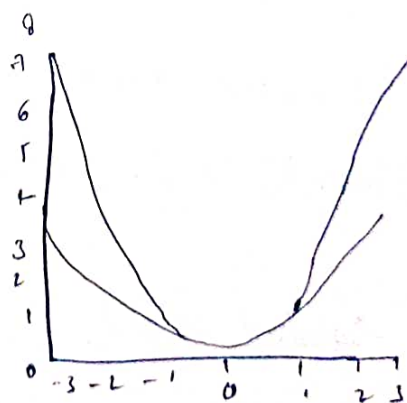
If there is outlier use MAE otherwise MSE

Advantage :-

- Used in robust regression.
- Outliers are handled properly.
- Local minima situation is handled here

Disadvantages :-

1. In order to maximize model accuracy, the hyperparameters δ will also need to be optimized which increases the training requirements.
2. It is complex.



$$L_{\delta}(a) = \frac{1}{2} a^2 \left(\delta \cdot |a| - \frac{1}{2} \delta \right)$$

for $|a| \leq \delta$, otherwise

$$L_{\delta}(y, f(x)) = \frac{1}{2} (y - f(x))^2$$

for $|y - f(x)| \leq \delta$, otherwise

RMSE

Unit

Output

Differentiable

$$\sqrt{\text{MSE}} = \sqrt{\frac{\sum_{i=1}^N \|y_i - \hat{y}_i\|^2}{N}}$$

Advantages

- It is helpful to ^{have} a single number ~~readable~~ to judge a model's performance, whether it be during training, cross-validation, or testing of the deployment.

Disadvantages:

Performance Metrics

1. R squared error
2. Adjusted R squared error
3. R Squared error

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



SS_{res} = Sum of Square Residuals

SS_{tot} = Sum of Square total

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

\bar{y} = Avg of y

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

$$R^2_{\text{Squared}} = 1 - \left\{ \frac{\text{Small number}}{\text{bigger num}} \right\} \rightarrow \text{Small number}$$

$$\boxed{\leq 1}$$

$\therefore 0.85 \Rightarrow 85\%$ accurate

$0.75 \Rightarrow 75\%$ accurate



$$R^2 = -ve$$

Performance of the model that you have created

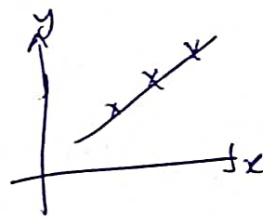
Can R^2 value be $-ve$?

\therefore

Model is very bad.

$$\boxed{R^2 = 1}$$

$\bar{y} = \text{Avg. of } y$



2. Adjusted R^2 squared :

Size of the house, City, Different location, No. of bedrooms, Gender, Salary, Price

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

$N = \text{No. of data points}$

$P = \text{No. of independent feature}$

$$R^2 \quad 65\%$$

$$R^2 \quad 75\%$$

$$R^2 \quad 88\% \quad \} - \text{No.}$$

$$R^2 \quad 90\%$$

$$\text{Adjusted } R^2 = 83\%, P = 1$$

$$= 73\%, P = 2$$

$$\text{Adjusted } R^2 = 85\%, P = n$$

Train :- very good accuracy [90%] [Low Bias]

Test :- Bad accuracy (50%) [High Variance]

→ Hyperparameter tuning.

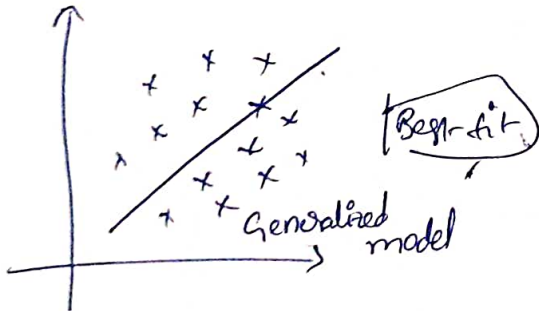
→ increase the data

Overfitting

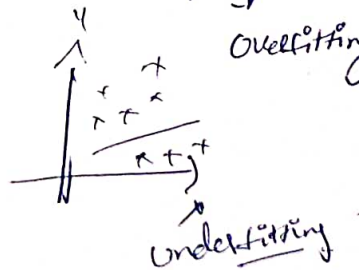
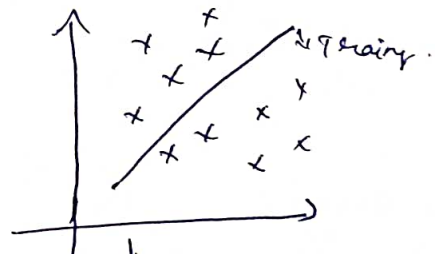
Train = Model accuracy is low [High bias]

Test = Model accuracy is low/high [low or high variance]

↓
model is Underfitting



Training data



• Linear Regression :-

1. Calculation = $J = mx + c$
2. Loss [MAE, MSE, RMSE, Huber loss]
3. Optimizing (m, c) → minimize the loss
(gradient decent)