

# Regression

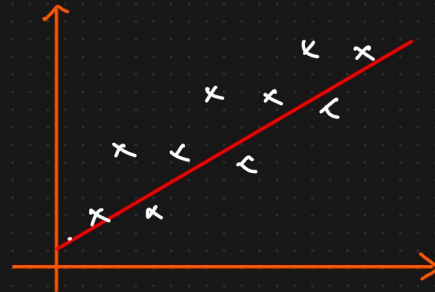
- ① Mean Square Error (MSE), MAE, RMSE ✓
- ② Ridge Regression ✓
- ③ Lasso Regression ✓
- ④ ElasticNet ✓
- ⑤ Practicals. [Simple Implementation]

Many regression models rely on distance metrics to determine the convergence to the best result. Even the definition of a "best" result needs to be explained quantitatively by some metric (COST FUNCTION).

Usually the metrics used are the Mean Average Error (MAE), the Mean Squared Error (MSE) or the Root Mean Squared Error (RMSE).

Distance calculated by taking difference of predicted and actual value should always be +ve. Distance obtained by difference can be +ve either:

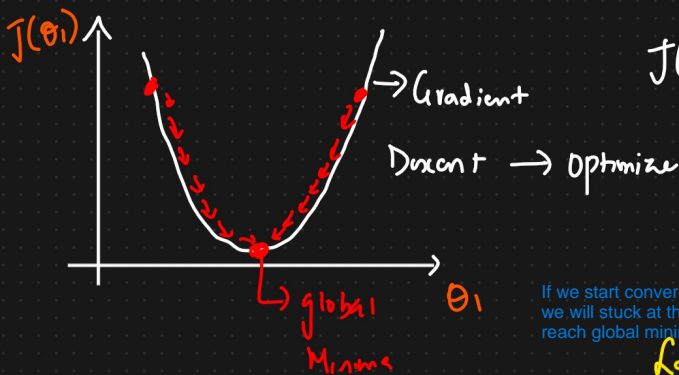
1. taking square=MSE
2. taking absolute difference=MAE
3. taking square root of squared quantity=RMSE



① Mean Squared Error (MSE)

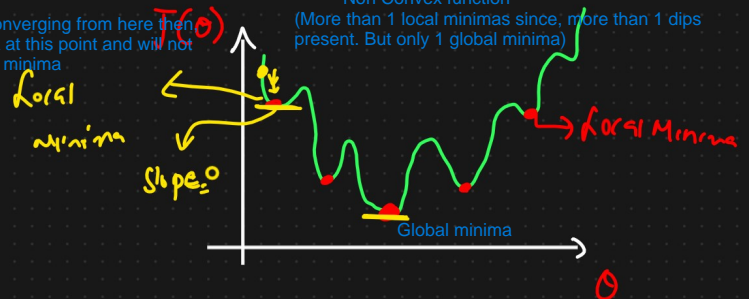
② MAE

③ RMSE



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

If we start converging from here then we will stuck at this point and will not reach global minima



① Mean Square Error [Cost function].

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

Parabolic function Convex function: only 1 local and global minima

$$ax + by + c = 0$$

$$by = -ax - c$$

$$y = \left[ \frac{-a}{b} \right] x + \left[ \frac{-c}{b} \right]$$

↓                      ↓  
mx + c

$$(a+b)^2 = a^2 + 2ab + b^2$$

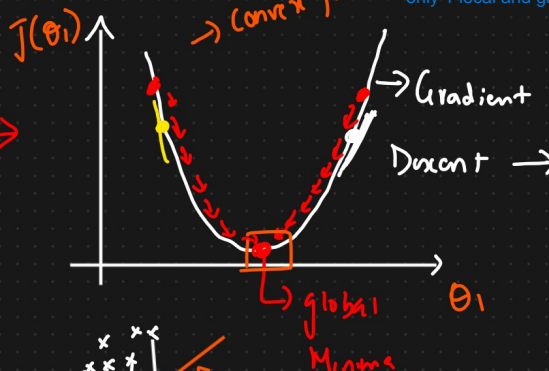
Below is representation of a straight line

$$ax + by + c = 0$$

$$y = mx + c$$

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

INR (INR)<sup>2</sup>  
[US] (US)<sup>2</sup>  
ERROR ↑↑  
MSE



$$MSE (y_i - \hat{y}_i)^2 \uparrow \uparrow$$

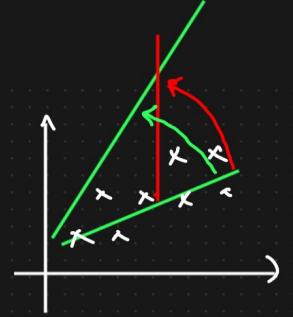
## Advantages

- ① It is differentiable.
- ② It has one local and one global minima.

Not robust to outliers means best fit line gets affected by outlier. Shift of best fit line is affected in huge way since we are squaring the difference. Ideally best fit line should move in small increments with varying data points but due to presence of outliers this shift increments drastically. Best fit line should shift in order to accommodate maximum data points.

## Disadvantages

- ① Not Robust to outliers  
*{Affected by outliers}*
- ② MSE changes the unit  
Since, in this we perform squaring operation. So, if unit is in cm for a feature it gets converted into cm square.



## ② Mean Absolute Error (MAE)

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

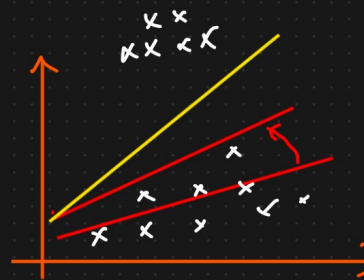
Advantages Since, here we are not squaring the differences, there will be shift in movement of best fit line in case outliers but this movement will be less as compared to MSE. **ERROR ↑**

- ① Robust to outliers

- ② It will be in same unit

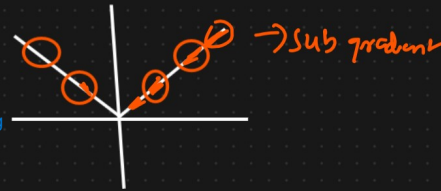
Due to above 2 advantages MAE is always preferred over MSE only in case there exists an outlier.

If there is no outlier then prefer MSE since it converges faster (since differentiable due to parabolic curve). On other hand MAE which will be straight line will converge at lower rate, thereby increasing the time complexity.



## Disadvantage

Time Complexity is more for optimization.



## ③ Root Mean Square Error

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \Rightarrow \text{Advantages And Disadvantages.}$$

Advantage:  
1. Unit does not change

Read this: <https://towardsdatascience.com/comparing-robustness-of-mae-mse-and-rmse-6d69da870828>

We can conclude that MSE and RMSE are more sensitive to outliers as compared to MAE.

Performance metrics =  $R^2$  and Adjusted  $R^2$

Cost function = MSE, MAE, RMSE, Huber Loss ← DL.

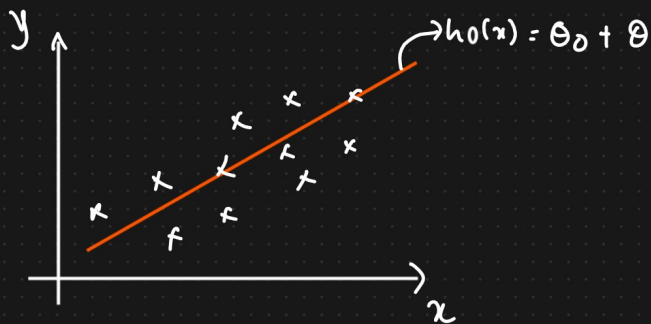
Huber loss cost function will be discussed in deep learning.



# Ridge, Lasso and ElasticNet

## Linear Regression

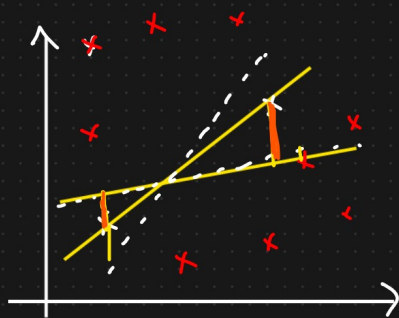
Simple regression already explained earlier



$$\text{Cost function} = \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x_i))^2$$

Ridge Reg follows L2 regularization used for reducing over fitting of data.

① Ridge Regression ( $L_2$  Regularization)  $\rightarrow$  Reducing Overfitting

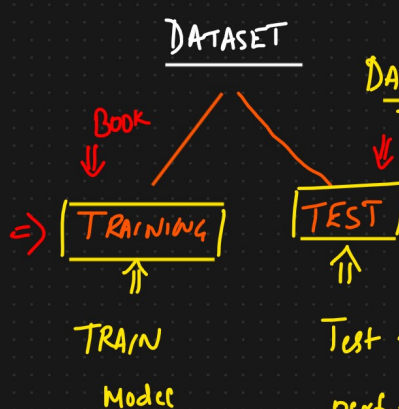


TRAINING  
Acc  $\Rightarrow 100\%$

Overfitting

Test Data  
Acc  $\downarrow \downarrow$

TRAINING  
DATA  $100\% \uparrow \uparrow$



DATA Leakage

Model should not have any info of Test data beforehand. Model should only know about Train data based on which model training is done. If model knows about test data as well then it would be a big blunder which is called data leakage.  
Eg: You know about questions coming in exams.

Test the Accuracy or performance of the model

Feature scaling should be done after train/test split. If done before then data leakage. Will be discussed in implementation.

In case of multiple regression this parameter will be squared sum of  $\theta_1, \theta_2, \dots$  so on.

Ridge Regression

$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x_i))^2 + \lambda \sum_{i=1}^n (\text{slope})^2$$

Zero

Hyperparameter

$$\lambda \geq 0$$

In general lambda is greater than zero. This value will decide the degree to which penalization happens which again will depend on problem statement.

$$\lambda = 1$$

Coefficients  
(Coefficient of x ie;  $\theta_1, \theta_2, \dots$  so on)

Important in case best fit line covers all training data points (overfitting) then also extra parameter is added that will make entire cost fun<sup>n</sup> non zero

$$+ 1 * [(\theta_1)^2 + (\theta_2)^2 + (\theta_3)^2] \Rightarrow \text{+ve value}$$

Over fitting occurs when accuracy of Train data increases whereas, accuracy of test data decreases (Reason self explanatory).

So, through ridge regression we are adding an extra parameter to cost function which is penalizing cost fun to avoid over fitting. When we add extra parameter then error increases making sure that best fit line does not entirely covers training dataset.

Cost fn = +ve value ↓↓↓

So to conclude after introduction of extra parameter accuracy of train data will decrease but model will become more capable to perform predictions based on any new varied data. That is test accuracy will increase. This is how we are reducing overfitting.

Coefficients ↓↓

Cost fun<sup>n</sup> here will always be +ve since in formula we are squaring the errors.

$\lambda = 1, 2, 3, 4, 10, 20, 30, 40, 50$

## Relationship between $\lambda$ and slope

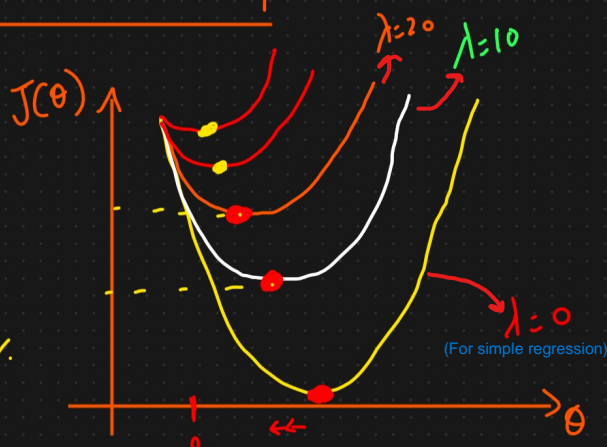
Observe that as we increase lambda, global minima is shifting left hand side and theta (on x axis) which is slope is getting shifted left or decreases.

Also on increasing lambda, error is also increasing which is represented by cost fun<sup>n</sup> on y axis.

TRAINING

Accuracy = 100%

TRAINING ERROR



$$Cost = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \cdot \sum_{i=1}^n (slope)^2$$

$\lambda \uparrow \uparrow$  Slope ↓↓

$\lambda \uparrow \uparrow \uparrow$  Slope ↓↓↓

$\lambda \uparrow \uparrow \uparrow$  ERROR ↑↑↑  
COST ↑↑↑

## ② Lasso Regression ( $L_1$ Regularization) → Feature Selection

( $L_1$  regularization)

multiplying a certain lambda that makes least correlated feature's coefficient zero and while doing so slopes of other features will decrease but will not become zero

$$Cost\ fn : \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n |slope| \leftarrow L_1\ Regularization$$

Eg:

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

$$h_{\theta}(x) = 0.52 + \boxed{0.65 x_1} + 0.72 x_2 + \boxed{0.12 x_3}$$

↓↓

Lasso Regression

Since  $x_3$  is least correlated feature it will get removed get removed by applying lasso regression. This is how lasso does feature selection by eliminating the least correlated feature.

unit change in  $x_1$

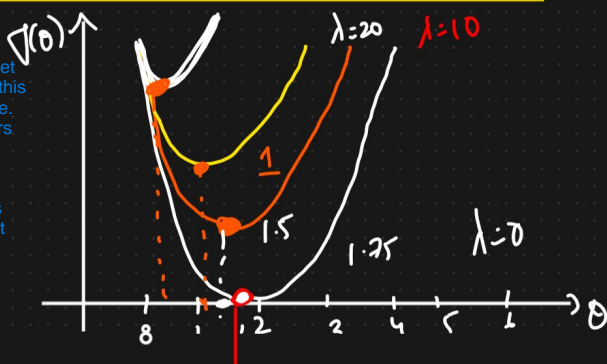
0.65 change in output

unit change in  $x_2$

0.12 change in output

## Relationship between $\lambda$ and |slope|

with increasing lamda the feature with least slope(theta value) will get eliminated as will tend to zero. In this way feature reduction is done here. With this other learning parameters will also decrease but will not become zero. The least important one will become zero and get eliminated from the list of features and in the process other important features will decrease but will not become zero.

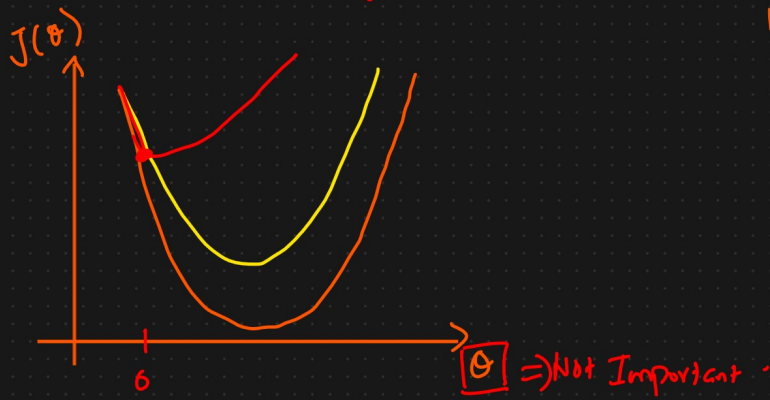


$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$h_{\theta}(x) = 0.52 + \boxed{0.65 x_1} + \boxed{0.72 x_2} + \cancel{\boxed{0.12 x_3}}$$

$\lambda = \dots$

$\boxed{0.12, 0.20, 0.50}$



$$h_{\theta}(x) = 0.52 + \boxed{0.65x_1} + \boxed{0.72x_2}$$

$\downarrow$                        $\downarrow$

$$= 0.52 + 0.45x_1 + 0.32x_2 + \boxed{\phantom{0.05x_1}}$$

③ Klassiker

$\left. \begin{array}{l} \rightarrow \text{① Reducing overfitting} \\ \rightarrow \text{② Feature Selection} \end{array} \right\} \text{Ridge + Lasso}$

$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \boxed{\lambda_1 \sum_{i=1}^n (\text{slope}_i)^2} + \boxed{\lambda_2 \sum_{i=1}^n |\text{slope}_i|}$$

$\lambda_1, \lambda_2$  [Hyperparameters].

$\downarrow$   
Reducing  
overfitting

$\downarrow$   
Feature selection.

Tuning of these hyperparameters is done in such a way so that we built a learning model that reduces overfitting and eliminates less correlated features.

Hyper tuning, gradient etc. these all are done just to build the best ml model.

