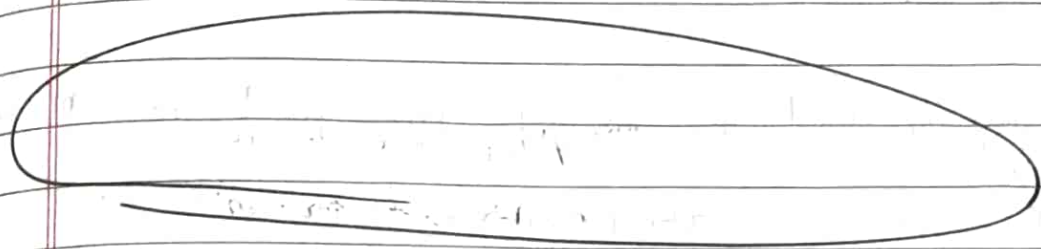


Simple linear regression → for continuous values

Supervised ML → Regression → for continuous values

Supervised ML → Classification → for categorical values.

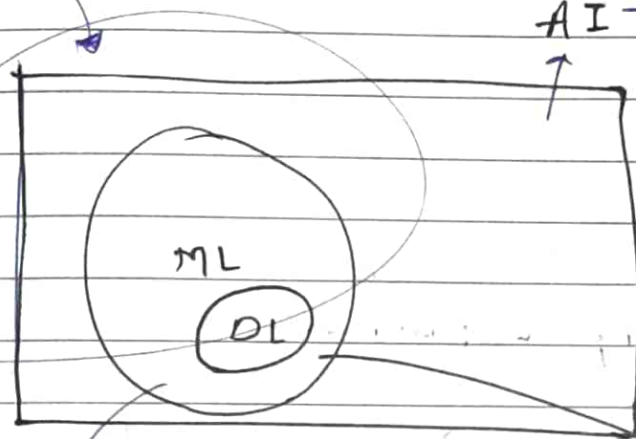


data science

AI → smart application that performs its own task without any human intervention.

- chatgpt

AI end product

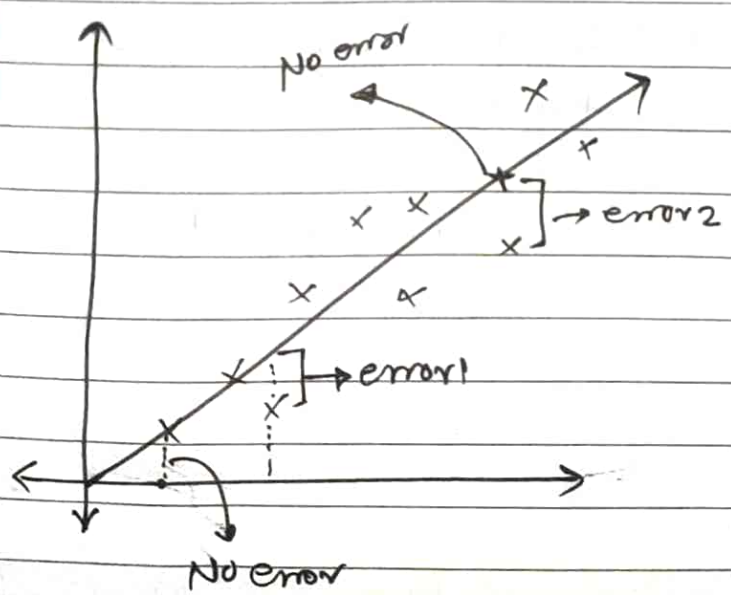


provide stats tool to perform explore, analyze, visualize and perform prediction.

Eg → Recommendation system

Netflix
AI powered application.

mimic human brain
• Multilayered Neural Network.



aim:-
create best fit line for data in such a way that error + error2 + ... errorn → have to be minimum,
• In case of two dimensional → it will be best fit line

○ for 4D there will be hyperplane.

equation of best fit line $\hat{y} = a_0 + a_1 x_i$ \rightarrow slope

$$\hat{y} = h_0(x) = a_0 + a_1 x_i$$

\downarrow
y-intercept

- By initializing a_0 and a_1 , we get the first fit plane.
- Then over the iteration we optimize it and for with line there is low error this line/ plane is best plane or best fit plane.

Cost function \rightarrow

error:

$$J(a_0, a_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x_i))^2 \quad \text{--- (mean squared error)}$$

\downarrow
predicted point

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

\downarrow \downarrow
no of data point actual data point predicted point

loss function \rightarrow

whenever we calculate the error for one point this is known as loss function

$$= (y_i - \hat{y}_i)^2 \quad \text{--- (1 data point)}$$

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

↓ To minimize this equation [value]
error

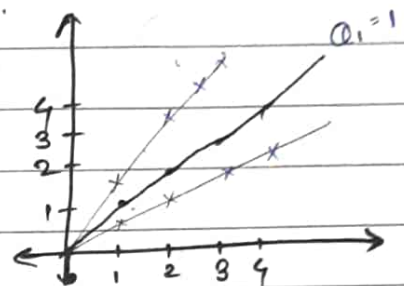
optimization:- [minimize the cost function].

$\theta_1 = 1$

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

$$= \frac{1}{4} [(1-1)^2 + (2-2)^2 + (3-3)^2 + (4-4)^2]$$

$= 0$ → minimum error
[can't be negative]



$\theta_1 = 0.5$

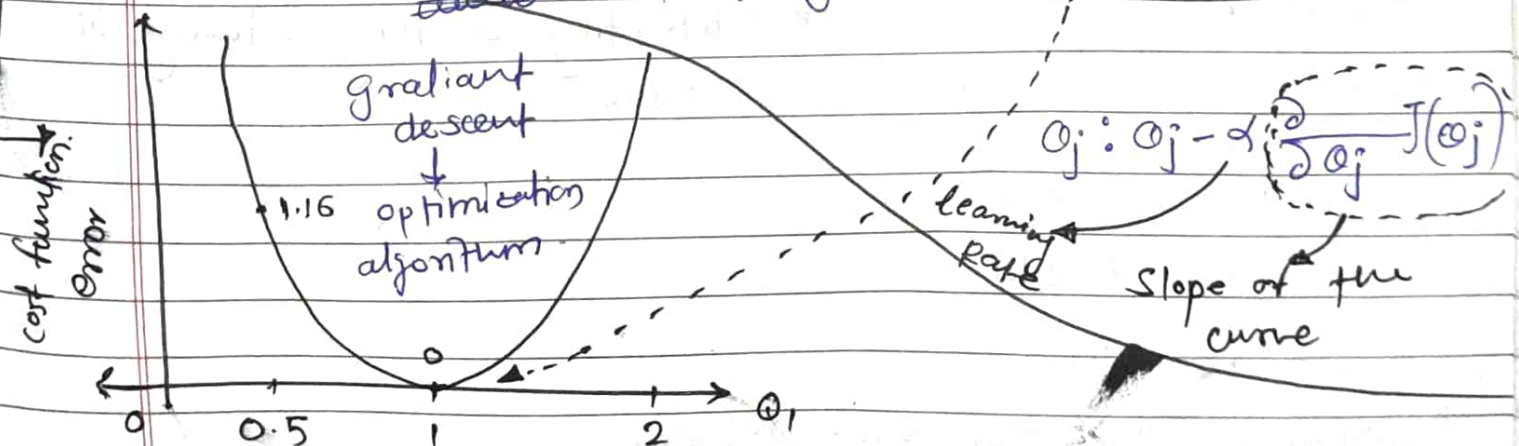
$$J(\theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

$$= \frac{1}{4} [(1-0.5)^2 + (2-1)^2 + (3-(0.5)(3))^2 + (4-2)^2]$$

$= 1.16$

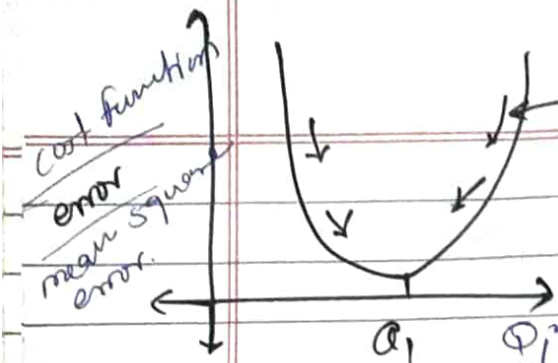
* Convergence algorithm:-

Repeat until ∇ Convergence global minima
 → Repeat the above process with different value of ' θ_1 ' till we getting to the minimum error point



→ at the value of $\theta_1 = 1$ we have minima point

descent gradient



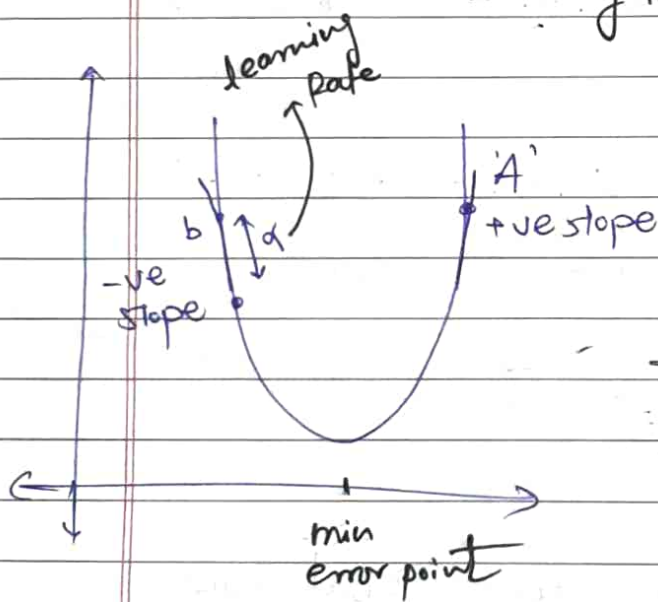
Travelling to minimize the error

→ optimization is used to minimize the error algorithm.

∴ The above travelling to make the error minimum is done by the optimization algorithm
Convergence algorithms

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

learning rate slope at point



for point 'A'

- To Reduce error we have to decrease ' θ_j '

with above equation

$$\theta_j := \theta_j - \alpha (\text{+ve value})$$

↓
slope

$$: \theta_j - \alpha (\text{+ve value})$$

→ Subtracting something Reducing θ_j value

for point 'b'

- To Reduce error we have to Increase ' θ_j '

with above equation

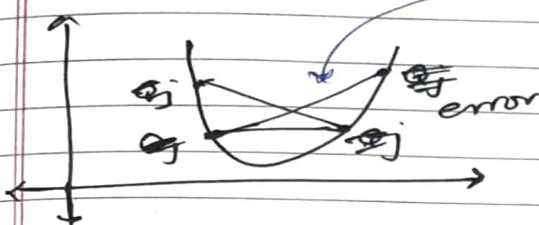
$$\theta_j := \theta_j - \alpha (-\text{ve slope})$$

$$\theta_j := \theta_j + \alpha (\text{slope})$$

→ adding something Increasing θ_j value

∴ If our learning Rate is too high
 α is too high

Then



• we can't able to reach at
 the minimum error point

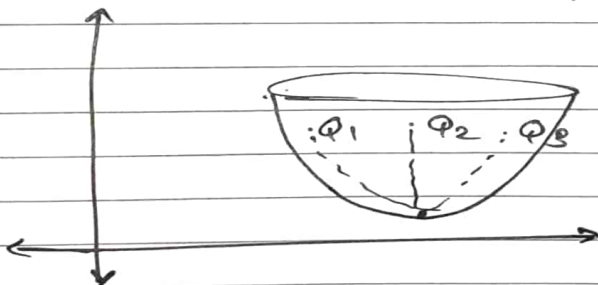
* if α is too small \rightarrow

Then it will take too much time to reach the minimum error point

* for multiple independent variable :-

$$h_{\theta}(x) = \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 \dots \theta_n f_n$$

-- multiple linear regression.



gradient descent curve.

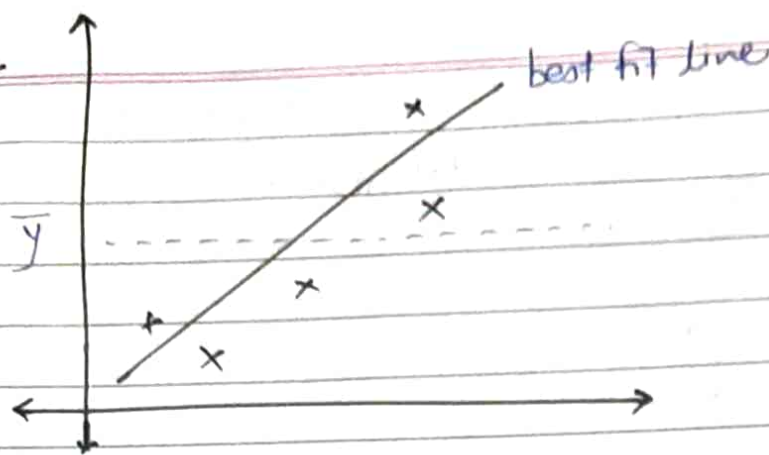
* Performance Matrix :-

classmate

Date

Page

* R-squared



$$R\text{-squared} = \frac{SS_{\text{regression}}(\text{Best fit line})}{SS_{\text{total}}(\text{average of } y)}$$

SS \rightarrow Sum of Square Residual or error

SS total \rightarrow Sum of Square total

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

This term will always
greater than numerator

$$\therefore R\text{squared} = 1 - (\text{value} < 1)$$

* Adjusted R-squared :-

$$\text{Adjusted R-squared} < R\text{squared}$$

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

N = no. of data points

p = No. of independent predictors

let's understand the adjusted $R^2 \rightarrow$ with example.

① let $R^2 = 80\%$, $N = 11$, $P = 2$

$$\text{Adjusted } R\text{-Squared} = 1 - \frac{(1 - 0.8)(10)}{11 - 2}$$

$$= 0.75$$

$$= 75\%$$

② Now, $P = 3$

~~Thus~~ The feature added is highly
Co-related then

Highly
Co-related

feature

R^2

adjusted R^2

2

80%

75%

3

85%

78%

4

86%

77%

$$R^2 = 85\%$$

$$\text{adjusted } R^2 = 78\%$$

③ Now, $P = 4$

- added feature not highly correlated
then

$$R^2 = 86\%$$

$$\text{adjusted } R^2 = 77\%$$

In Linear Regression we use Mean square error as loss function.

classmate

Date _____
Page _____

$$MSE = \frac{1}{n} \sum_{i=1}^n [y_i - \hat{y}_i]^2$$

advantage of MSE

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

Disadvantage of MSE

- ① Not robust to outliers

[to reduce the error occurred by the outlier]

- algorithm will try to reduce it change the bestfit line]

- ② MSE changes the unit

[by making the square of it].

② → Mean Absolute Error ⇒ (MAE)

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

→ Advantage of this ⇒

- ① Robust to outliers
- ② It will be in same unit.

disadvantages:

- time complexity is more for optimization
- instead of gradient descent you get



③ Root mean square error ⇒

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

advantage

disadvantage

* Ridge, Lasso and elasticnet :-

Cost function for linear Regression

$$= \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x))^2$$

Note :-

always do feature scaling after the train test split only

① Ridge Regression → (L2 Regularization)

- we are basically overfitting,
 ↓
 Reduce

Cost function :-

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

for parameter 2

Hyperparameter

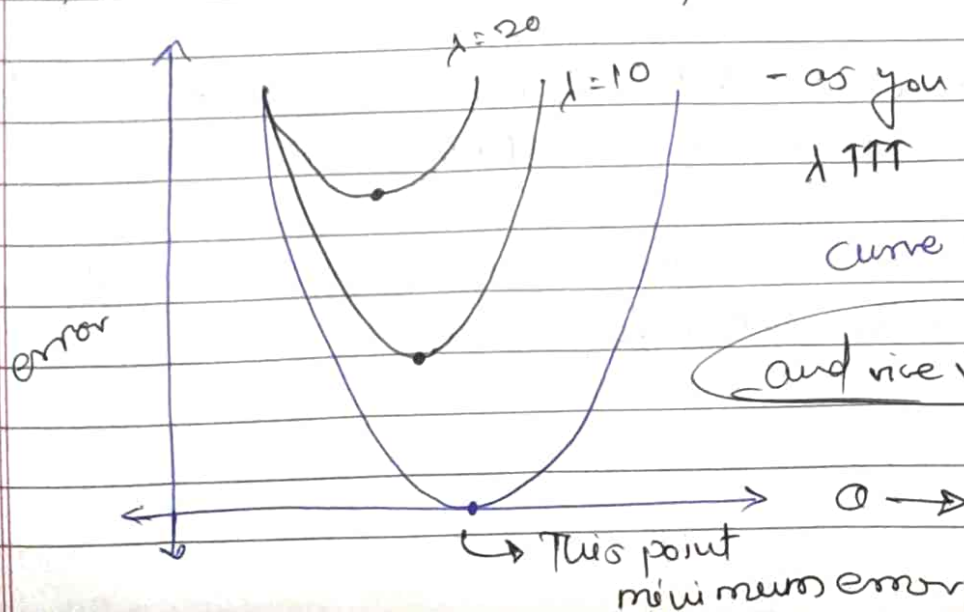
$$[(\theta_1)^2 + (\theta_2)^2 + (\theta_3)^2]^2$$

for parameter 1 for parameter 3

overfitting

- due to this we always adding error in the Cost function of linear Regression
- So basically we reducing the overfitting.

* Relation between 'λ' and slope :-



- as you can see

λ ↑↑↑ slope ↓↓

curve

error ↑↑↑

and vice versa

This point minimum error

② Lasso Regression \Rightarrow (L1 Regularization)

- Mainly used for feature selection

classmate

Date _____
Page _____

Cost function \Rightarrow

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

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$h_0(x) = 0.52 + 0.65x_1 + 0.72x_2 + 0.12x_3$$

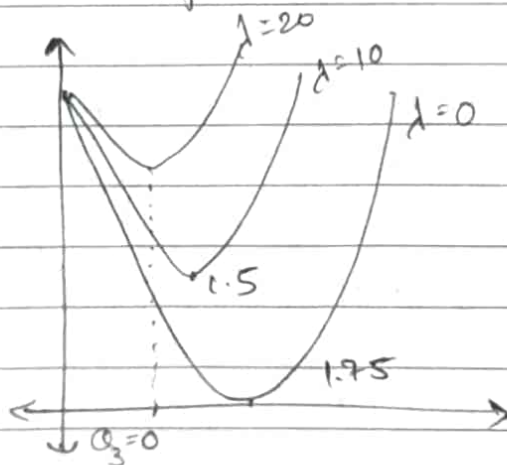
x_1
 x_2
 x_3

} features

with unit change in x_i there is 0.65 change in output (prediction)

* So the x_3 contribute much less to predict the output

• Relationship between λ and $|\text{slope}|$



will reach to '0' because we use $||$ absolute value not like Ridge.

• So it will remove the impact of less correlated feature internally

* ElasticNet \Rightarrow

- Combination of both

\rightarrow Ridge \rightarrow Reducing overfitting
 \rightarrow Lasso \rightarrow feature selection

classmate
Date _____
Page _____
Reducing overfitting.

Cost function $\Rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x)_i)^2 + \lambda_1 \left(\sum_{j=1}^n (\text{slope})^2 \right) +$

$+ \lambda_2 \sum_{j=1}^n |\text{slope}|$

\downarrow
feature selection

* sns.pairplot(df)

- plot the scattered plot for all feature in pair.