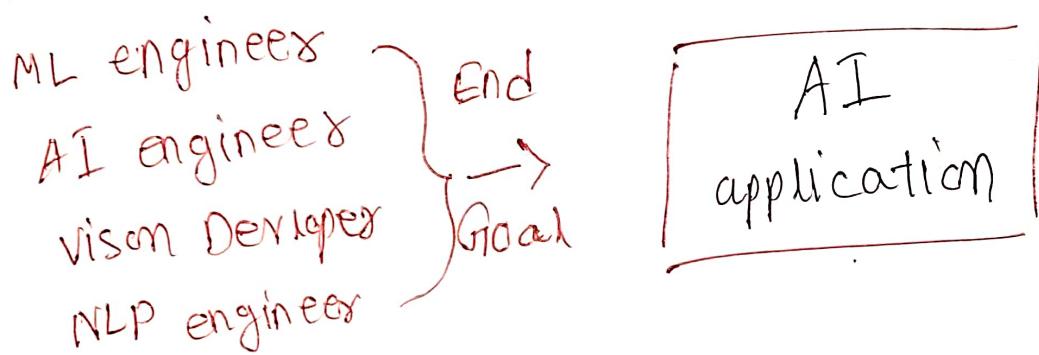


Introduction To Machine Learning

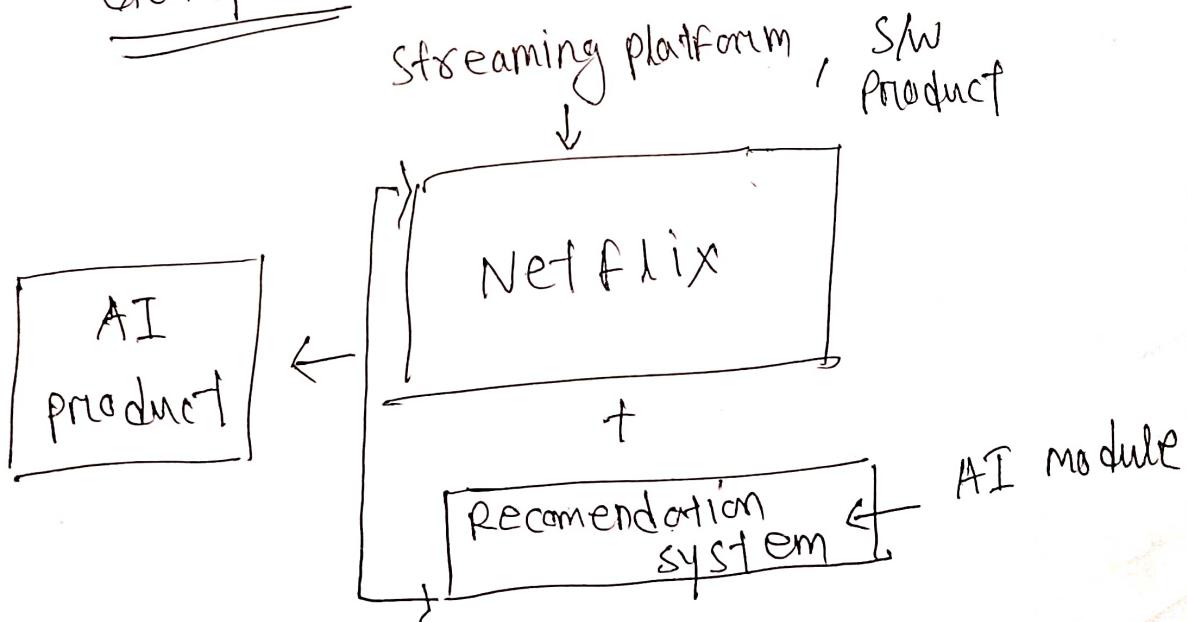
Agenda

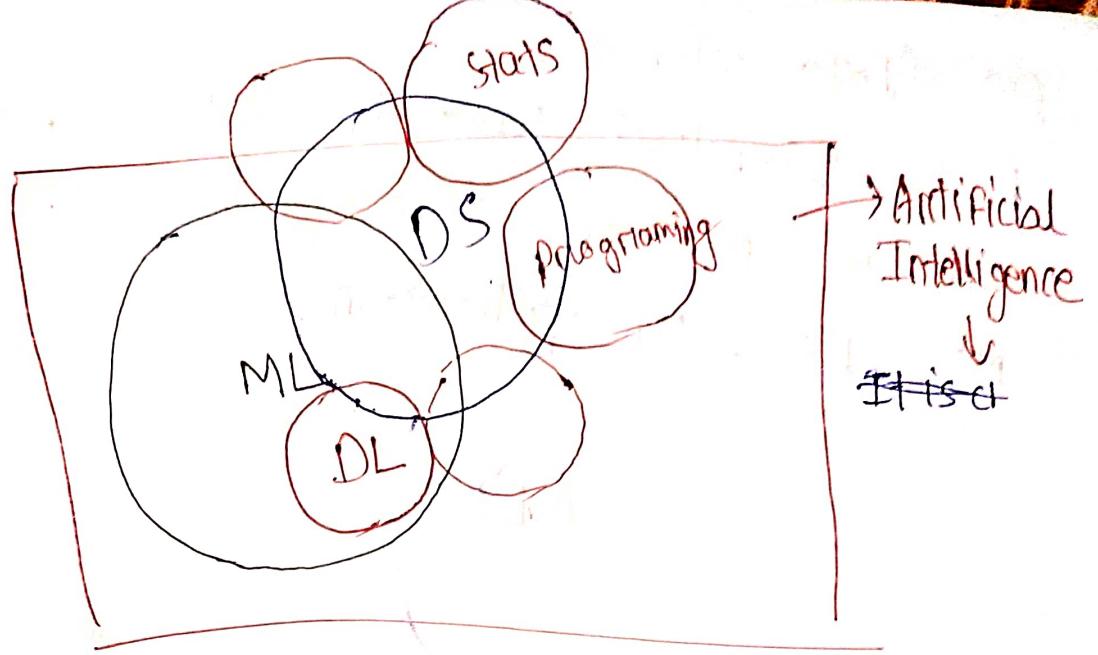
- ① Machine Learning Introduction
- ② AI vs ML vs DL vs DS
- ③ Simple Linear Regression →
Mathematical Intuition

AI Product



example





AI (Artificial Intelligence)

It is creating an application where it performs all its task without any human intervention.

[Human try to make this product better and better]

examples

AI chatbots, Alexa, self-driving car, YouTube

ML (Machine Learning)

ML provides stats tools to analyze, visualize perform predictions and other task with the help of Data.

DL (Deep Learning)

Around 1950's researchers thought's that we can make machine learn, how we are able to learn.

DS (Data science)

Data scientist has to deliver what is needed, it may be related to ML, DL, AI application.

ML
T

Supervised

Regression

① Linear Regression

② Polynomial "

③ SVR

④ Decision tree Regressor

⑤ Random Forest Regressor

⑥ Xgboost

⑦ Naive Bayes

Clustering Algorithms

DBScan, k-means

hierarchical clustering
(grouping)

① Logistic Regression

② SVM

③ Decision tree classifier

④ Random forest classifier

Supervised

O/P will be given

make classification,

Unsupervised

O/P not given

Clusters(groups)

Regression

Supervised

Degree

exp

B.E

7

P.H.D

2

-

-

-

-

-

Input on Independent features

Number of study hours

9

7

2

1

No. of play hours

1

3

5

6

(O/P) \leftarrow continuous feature

Salary

50K

70K

-

-

-

-

-

-

-

-

-

-

-

-

-

-

Regression problem

Dependent feature.

O/P

Person pass/fail

pass(1)

1

0 (fail)

0

Classification problem

- Flight price prediction (Regression)
- Tomorrow going to Rain or not (Classification)
- Air quality index (Prediction or Regression)

Simple Linear Regression

one independent \rightarrow one dependent feature

examples

No. of rooms \rightarrow Price

Model \leftarrow Years of experience & Salary

Predict \rightarrow Salary Based on 1/p years.

Years of Exp

Salary

30K

30K

20OK

-

-

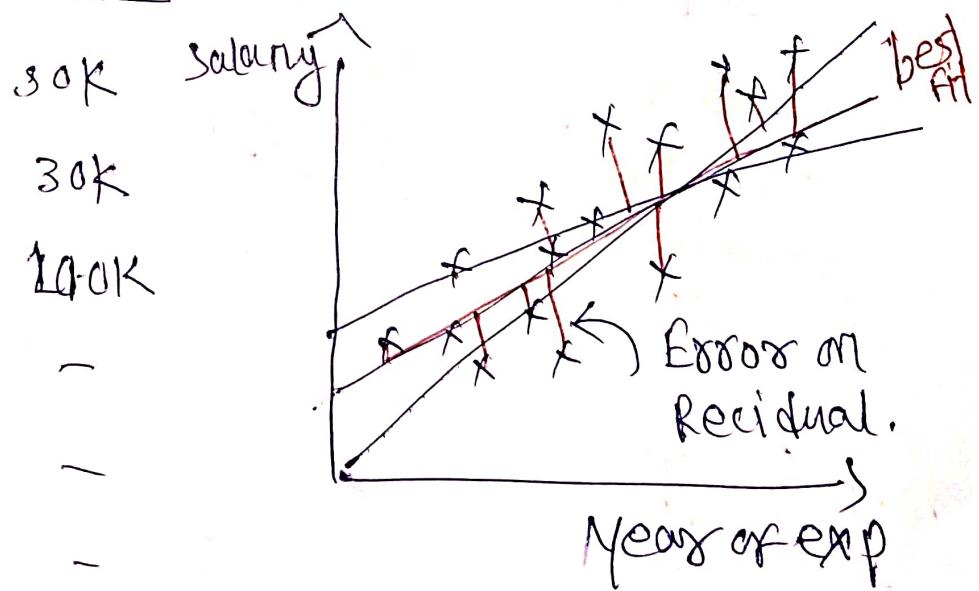
-

-

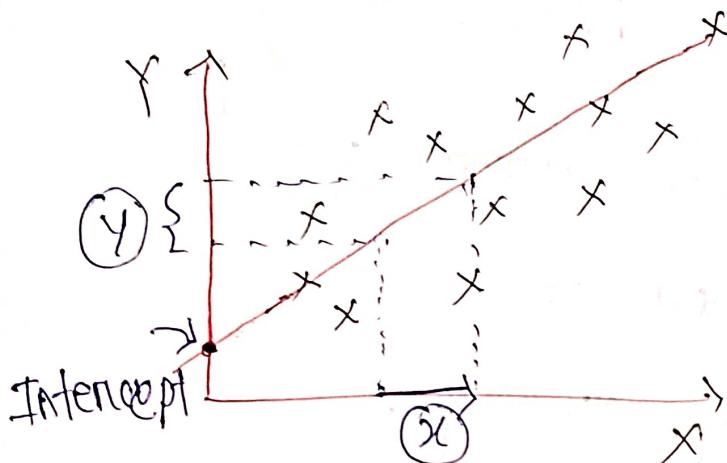
-

-

-



In simple linear regression we try to find the best fit line such that error (distance between all points and prediction (actual)) is minimum.



Equation of straight line

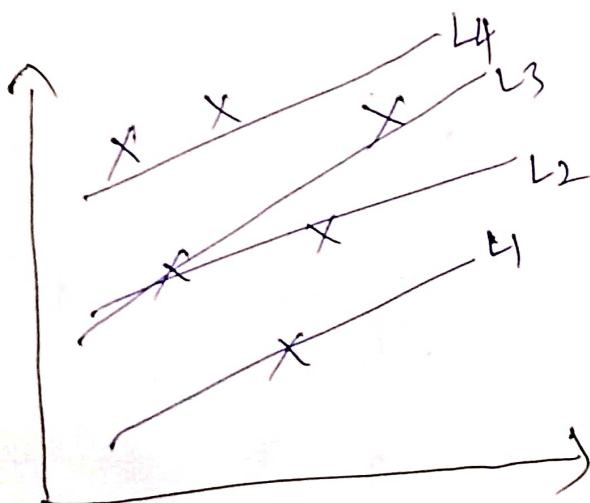
$$Y = Mx + C$$

$$y = B_0 + B_1 x$$

slope (unit movement in Y axis,
when unit movement
in x axis)

$$h_0(x) = B_0 + B_1 x$$

Intercept ($Y=2$, when $x=0$)



Training of
the model.

→ To find the best fit line we change in θ_0 and θ_1 and predict the best fit line.

→ Cost function

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

mean square error ↓↓↓

$J(\theta_0, \theta_1)$ → cost function

$\frac{1}{m}$ → for finding average.

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

Let

$$\theta_0 = 0$$

$$h_\theta(x) = \theta_1 x$$

when

$$\theta_1 = 1$$

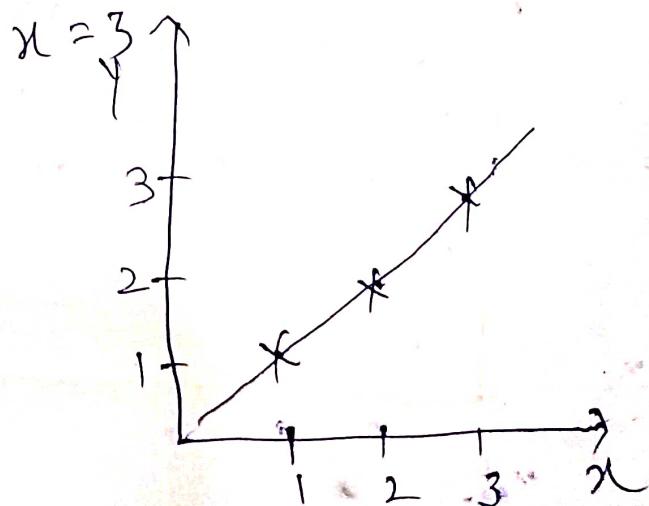
$$\therefore h_\theta(1) = 1 \times 1 = 1 \quad x=1$$

$$h_\theta(2) = 2 \quad x=2$$

$$h_\theta(3) = 3 \quad x=3$$

<u>Training Data</u>	
x	y
1	1
2	2
3	3

→ Both actual and predicted points are equal.

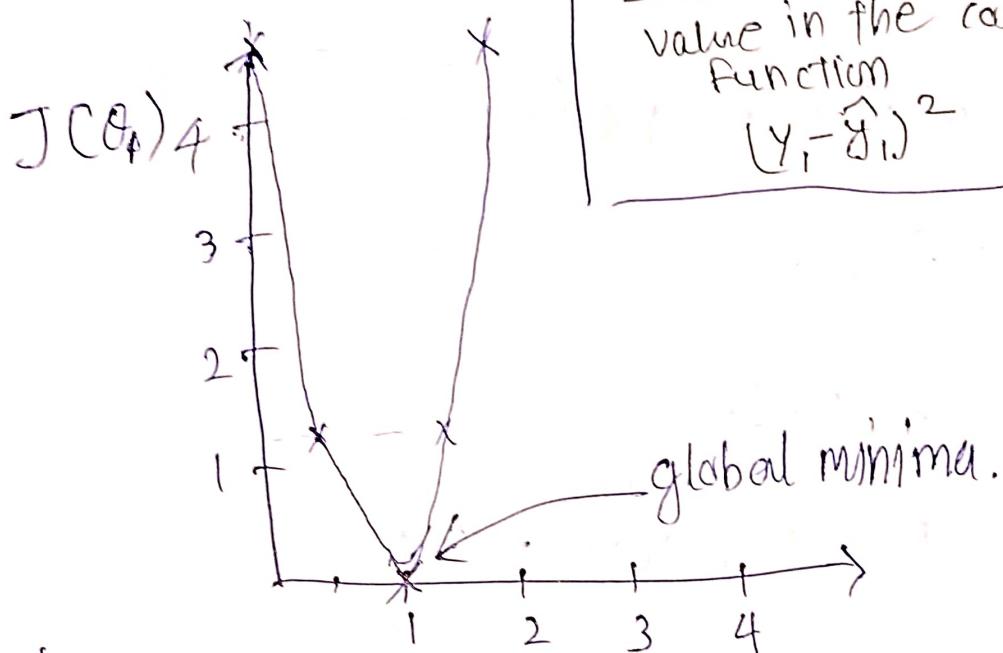


error

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

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

$$= 0$$



when

$$\text{set } \theta_1 = 0.5$$

$$h_\theta(x) = 0.5x \quad x=1$$

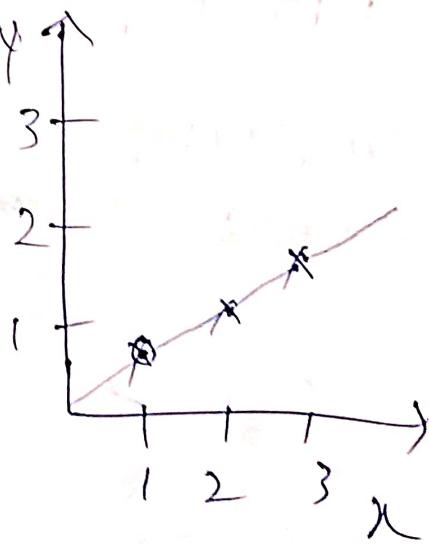
$$h_\theta(x) = 1 \quad , x=2$$

$$h_\theta(x) = 1.5 \quad , x=3$$

$$J(\theta_1) = \frac{1}{3} ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2)$$

$$= \frac{1}{3} (0.25 + 1 + 2.25)$$

$$= \frac{3.5}{3} = 1.16$$



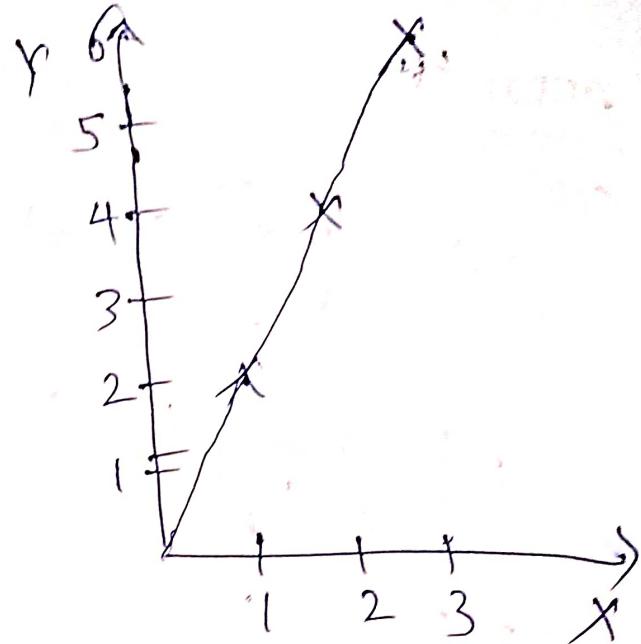
when

$$\theta_1 = 2$$

$$h_0(x) = 2, x=1$$

$$h_0(x) = 4, x=2$$

$$h_0(x) = 6, x=3$$



$J(\theta_1)$

$$= \frac{1}{3} ((2-1)^2 + (4-2)^2 + (6-3)^2)$$

$$= \frac{1}{3} (1^2 + 2^2 + 3^2)$$

$$= \frac{14}{3} \approx 4.66$$

when $\theta_1 = 0$, $J(\theta_1) \approx 4.66$ similarly.

At $\theta_1 = 1$, $J(\theta_1) = 0$ \downarrow minimum.

Like that we have to reach the minimum point called global minima.

→ The technique for reaching global minima is called convergence algorithm.

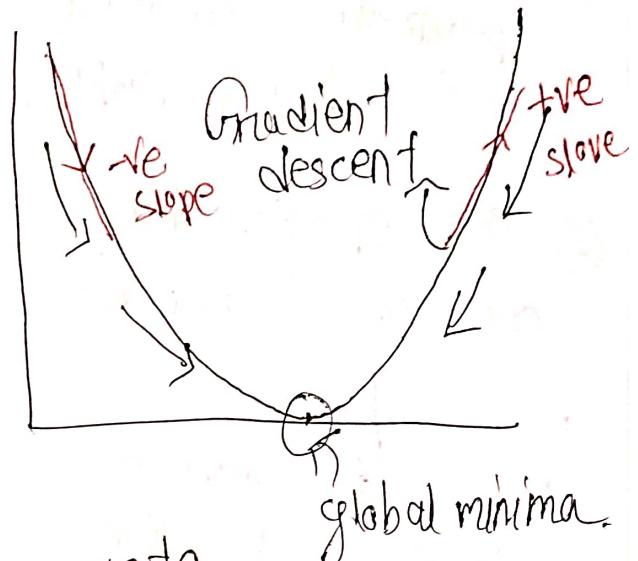
convergence Algorithm (optimize the changes of θ_1 value)

Algorithm

Repeat Until convergence

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

α = Learning rate



α / learning rate decides the speed to reach the global minima. It should not be too high or it moves in zig-zag speed and not able to reach global minima. It should not be too low also. generally $\alpha = 0.001$

Cost functions

① MSE

② MAE

③ RMSE

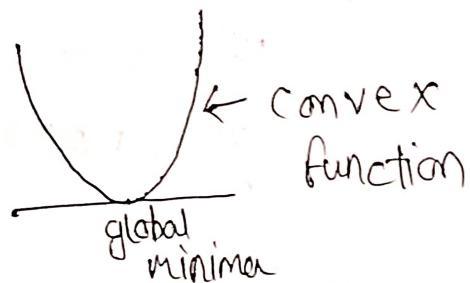
① MSE (mean square error)

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

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

Advantages of MSE

- ① As MSE is a quadratic equation it is differentiable.
- ② This equation has one global minima.



Disadvantage of MSE

- ① This is not robust to outliers.
- ② Due to quadratic equation unit of output feature changes or penalizing error.

MAE (Mean absolute error)

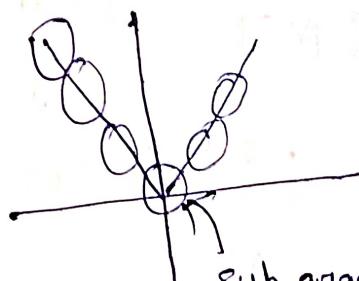
$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Advantages

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

Disadvantage

- ① convergence usually takes more time.
optimization is a complex task



because of sharp zero value.

RSME (Root Mean Square Error)

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

Advantages

- ① Unit will be same due to square root.
- ② It is differentiable

Disadvantages

- ① This is not robust to outliers.

Performance Matrix

- ① R Squared
- ② Adjusted R squared

- ① R squared

$$R^2 = 1 - \frac{SS_{\text{Res}}}{SS_{\text{Total}}}$$

SS_{Res} = Sum of square Residual.

SS_{Total} = Sum of square Average.

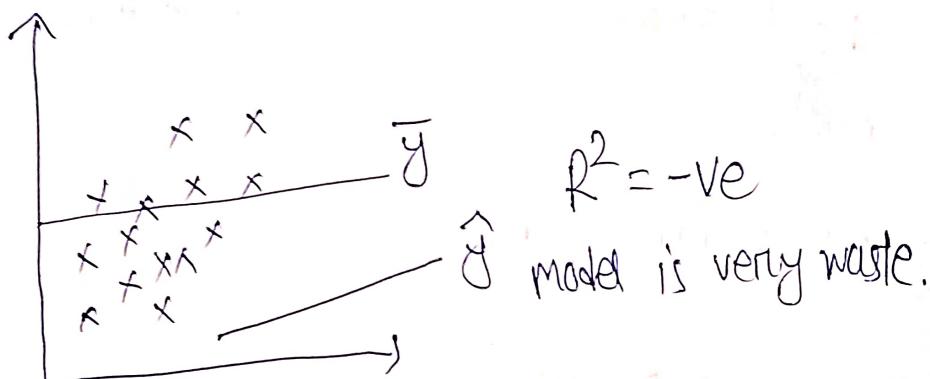
$$R^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

\bar{y} = Average of y

$$R^2 = 1 - \left(\frac{\text{small num}}{\text{large num}} \right) \} \text{ small numbers}$$

→ R square measures performance of the model.

Q. Can R^2 squared value -ve?



$$R^2 = 1 - \left(\frac{\text{large num.}}{\text{small num.}} \right) - \text{large } (C > 1)$$

$= (-ve)$ ✓

② Adjusted R squared

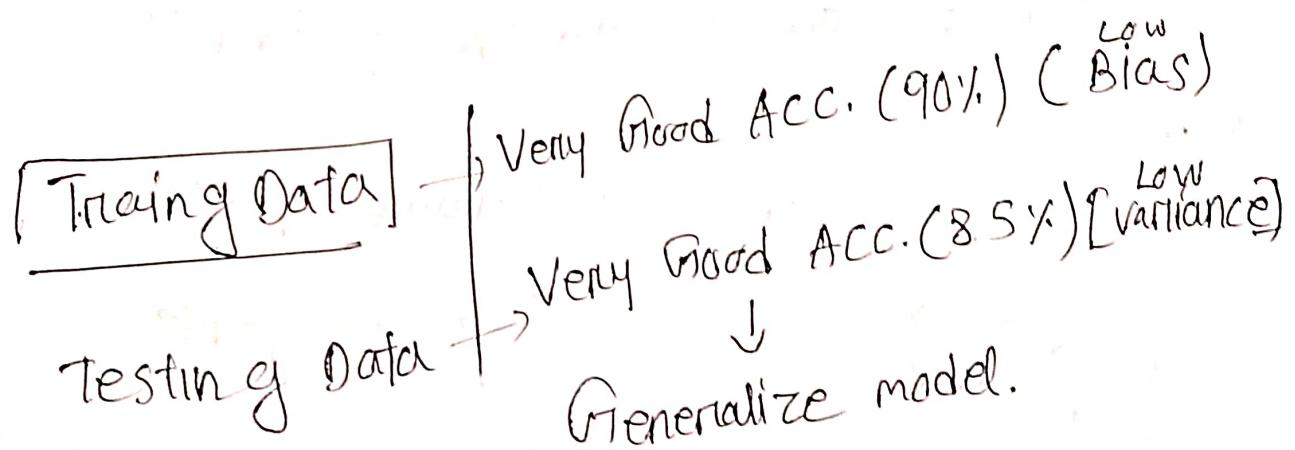
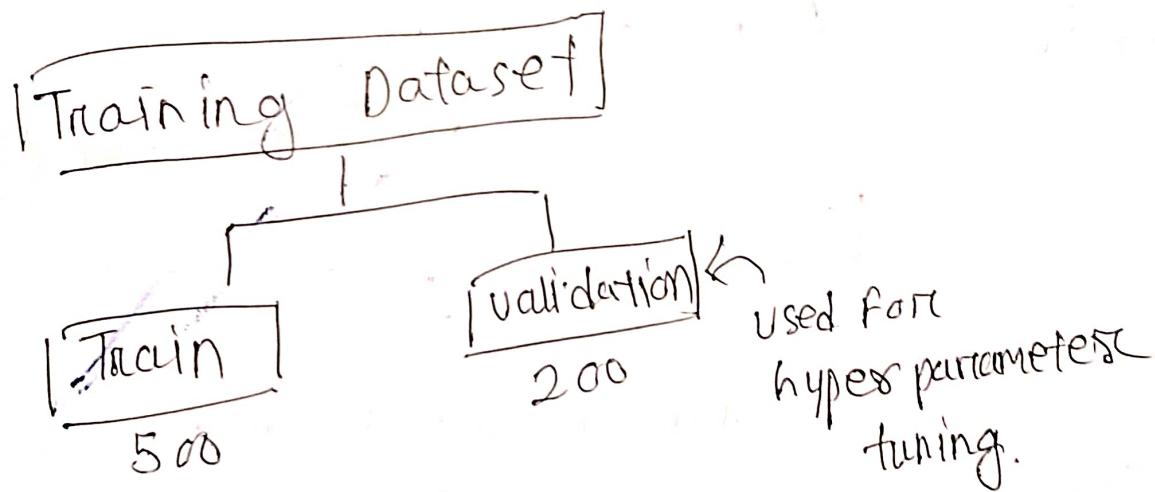
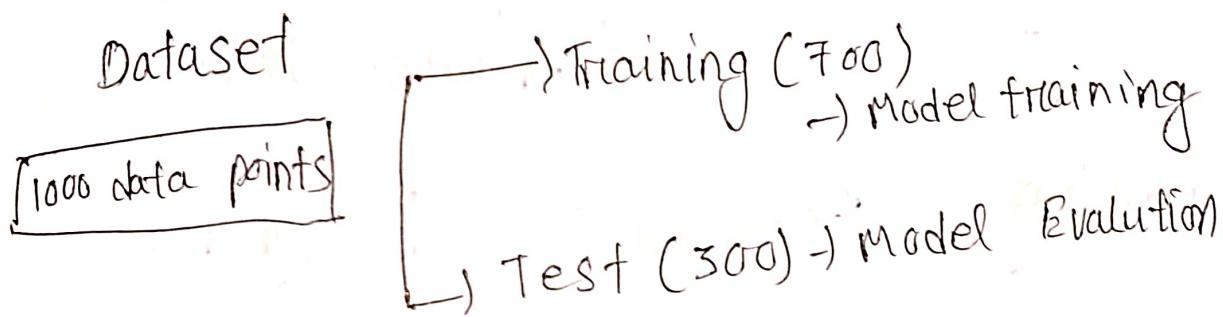
	I/P			O/P	
	size of house	city location	No. of bedrooms	Gender	Price
R^2	65%	75%	88%	90%	
Adjusted R^2	65%	75%	88%	85%	

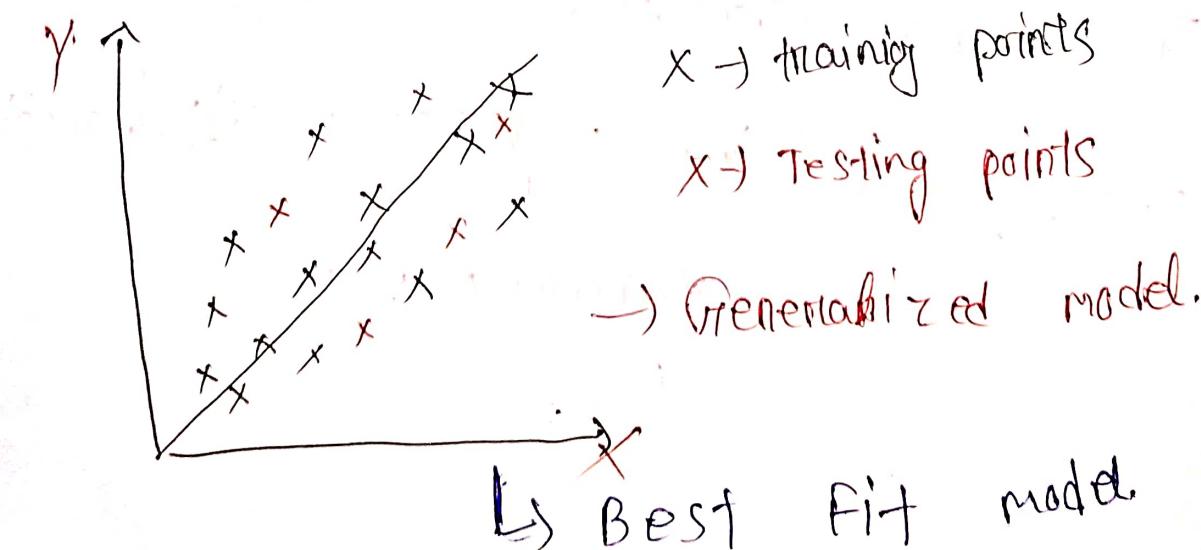
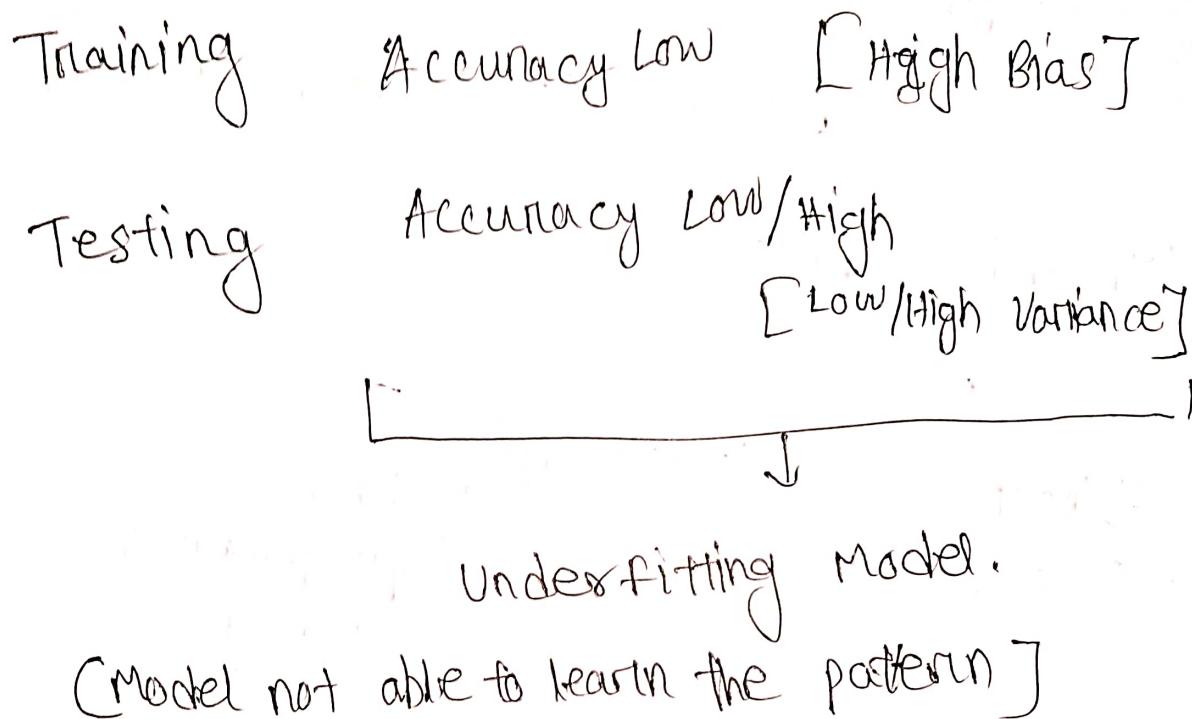
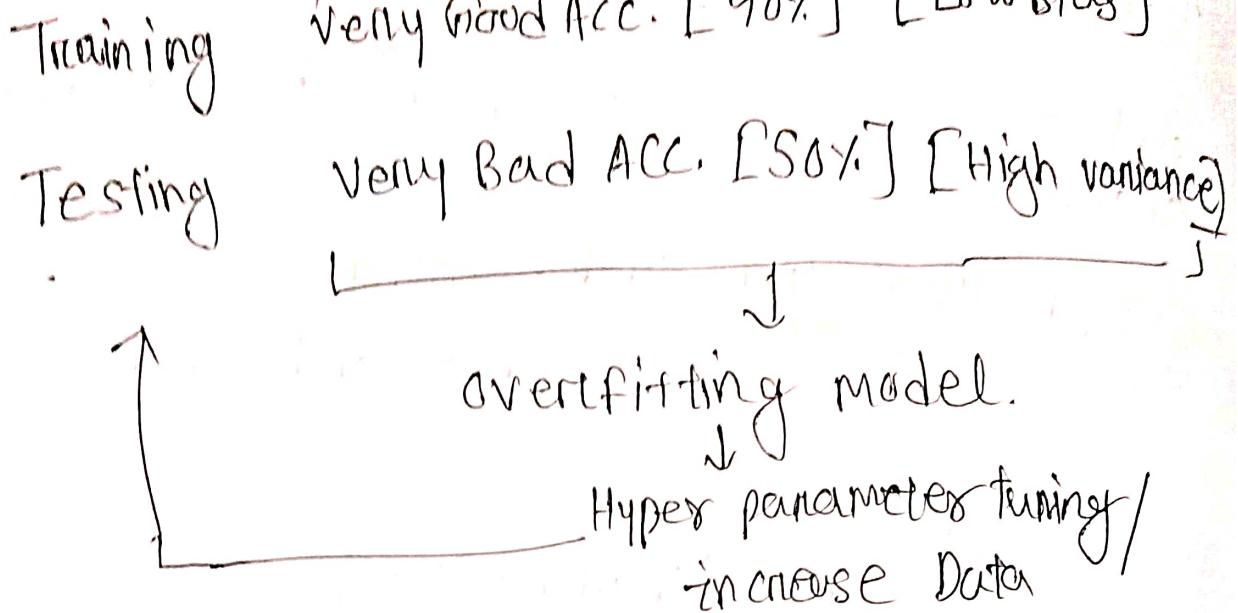
By adding a valuable columns accuracy (R^2) generally increases. with adding unnecessary column it increases slightly (R^2)

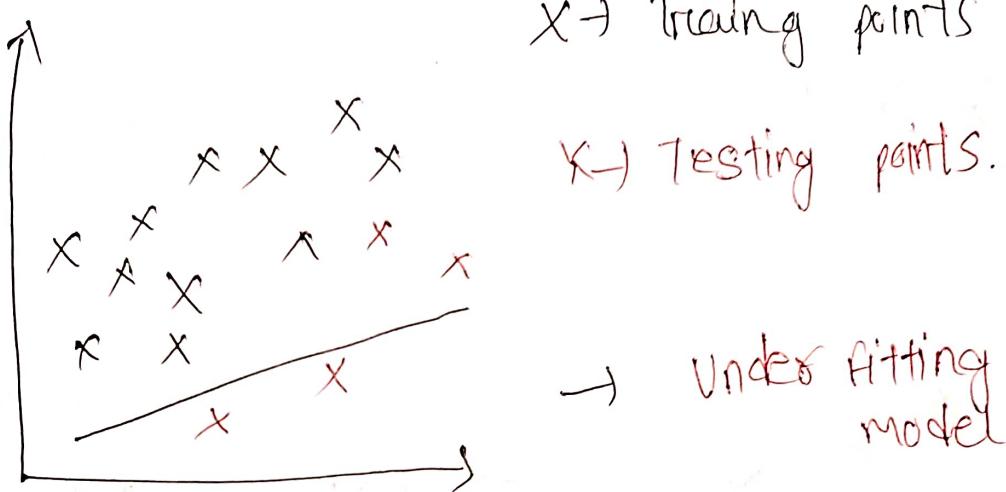
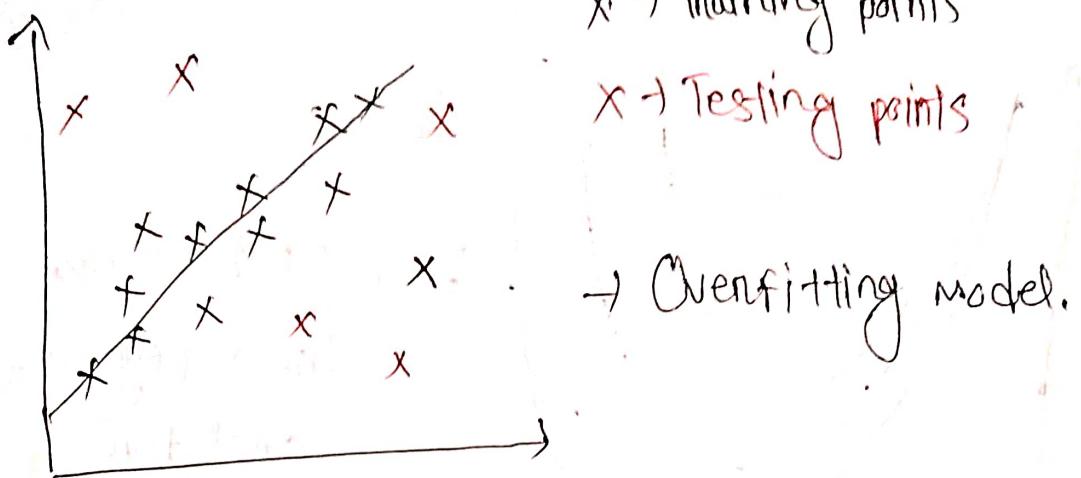
$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

But in case of adjusted R^2 it decreases slightly.

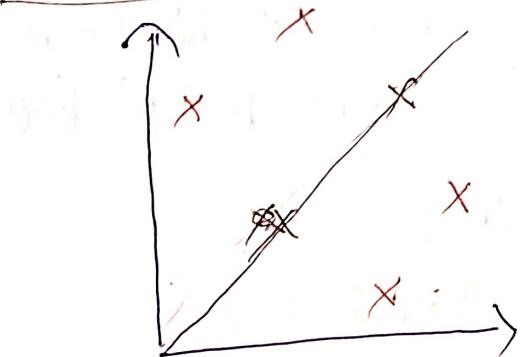
Overfitting and Underfitting
(Bias and Variance)







~~Ridge Regression (L₂ Regularization)~~



overfitting

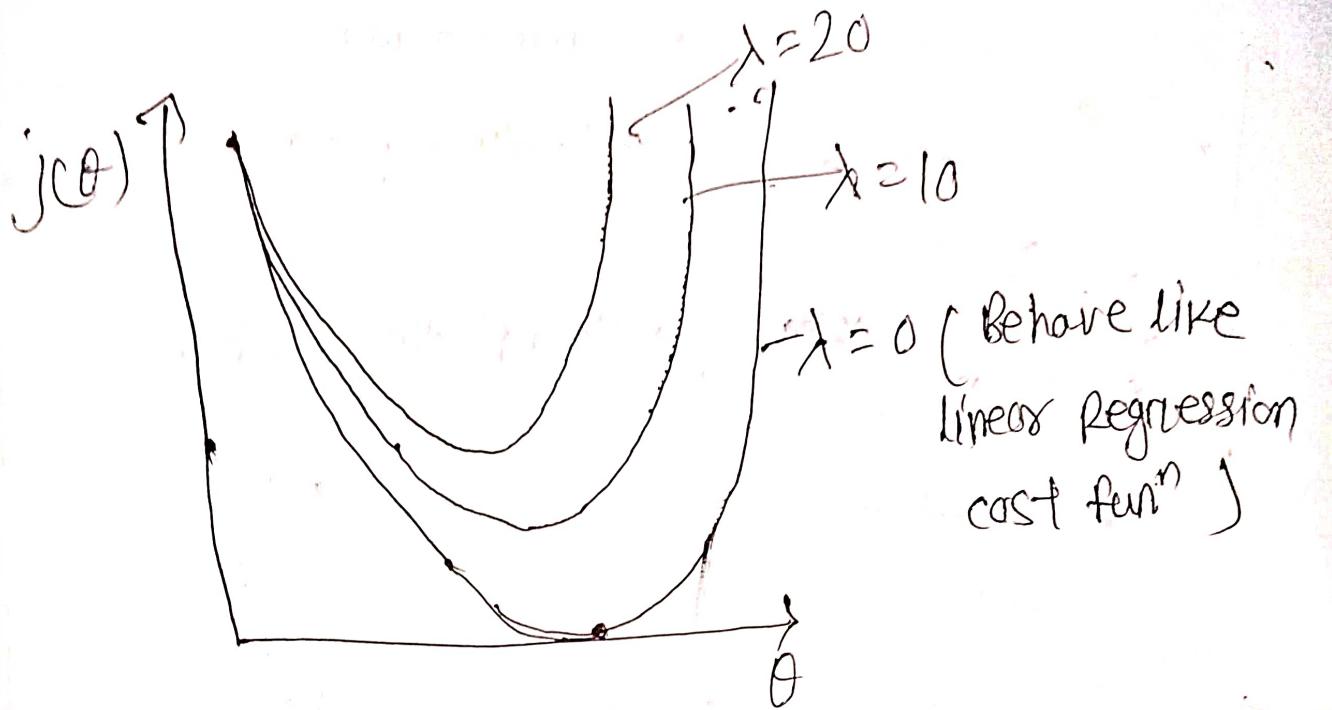
Training accuracy = 95%

Test accuracy = 50%

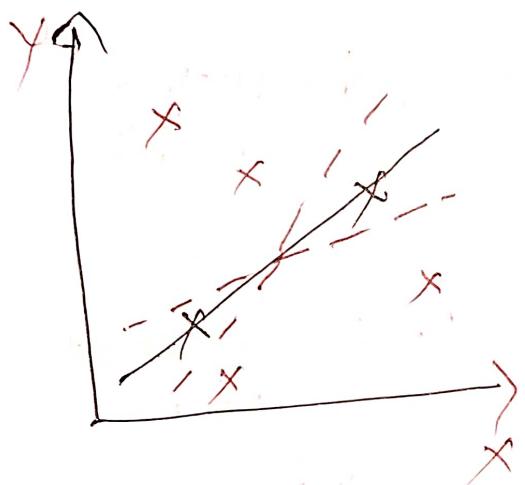
cost fn = $\frac{1}{m} \sum_{i=1}^n (h(x)^i - y^i)^2 + \lambda \sum_{j=1}^n (\text{slope})^2$

(~~linear reg.~~)

λ = hyperparameter



$$\lambda \propto \frac{1}{\theta}$$



$$\text{cost fun}^n = 0 + \lambda (\text{slope})^2$$

= +Ve value
to decrease this we need to modify best fit line.

→ Reduce overfitting
(Never Overfit)

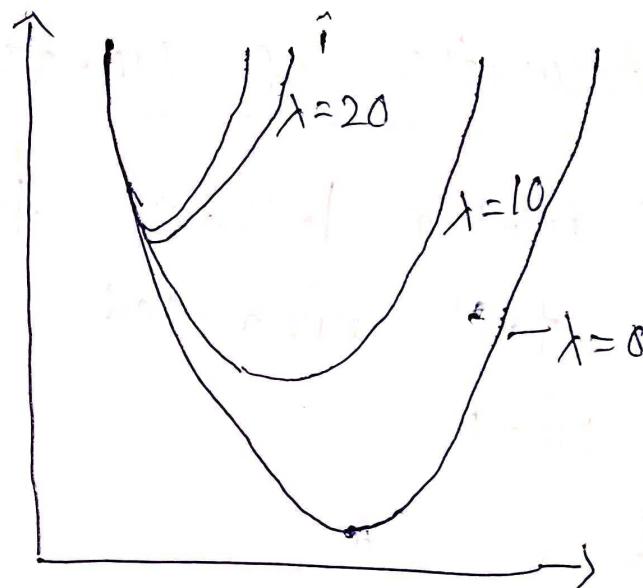
Note
 $\theta \neq 0$

LASSO Regression (L_1 regularization / L_1 Norm)

Used for feature selection.

$$\text{cost fun} = \frac{1}{m} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m |\text{slope}|$$

As λ increases
slopes decrease
so some features
deleted (not important
features)



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

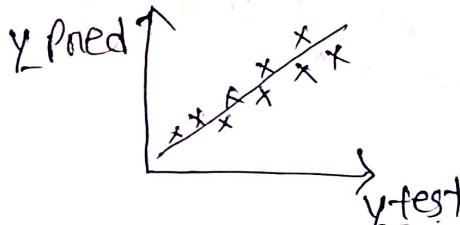
$$= \theta_0 + 0.54x_1 + 0.82x_2 + 0.10x_3$$

↓
feature removed.

* Elastic Net [combination of L_1 and L_2 Norm]

$$\text{cost fun}^n = \frac{1}{m} \sum_{i=1}^n (h_\theta(x_i) - y_i)^2 + \lambda \left(\underset{\text{Ridge}}{\sum_{i=1}^n (\text{slope})^2} + \underset{\text{Lasso}}{\lambda_2 \sum_{i=1}^n |\text{slope}|} \right)$$

Assumption of Linear Regression

- ① Relation between test truth and test predicted data should be linear.

- ② Residual = $y_{\text{test}} - y_{\text{pred}}$
→ Residual should follow approximate gaussian distribution.
- ③ Plot between prediction and residual follow uniform distribution.