

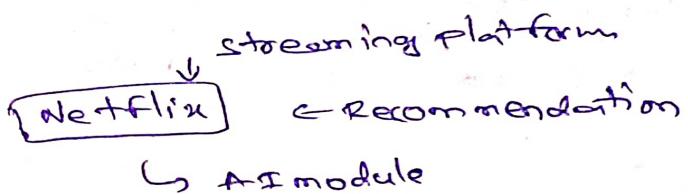
Introduction to machine learning

Agenda

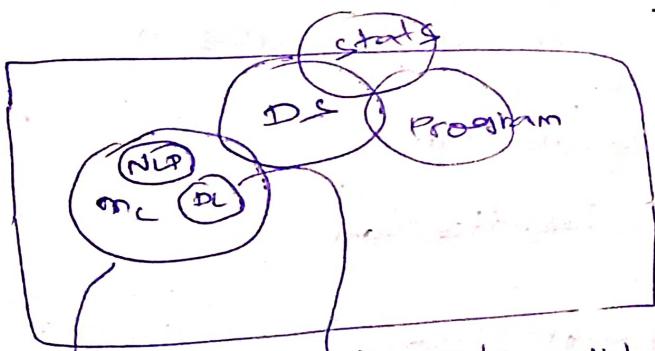
- ① machine learning introduction
- ② AI vs ML vs DL vs DS
- ③ simple linear regression — mathematical Introduction

AI product

Netflix: Now a day we all are seeing Netflix when we see some movies it will record some movies is called AI product



AI vs ML vs DL vs DS



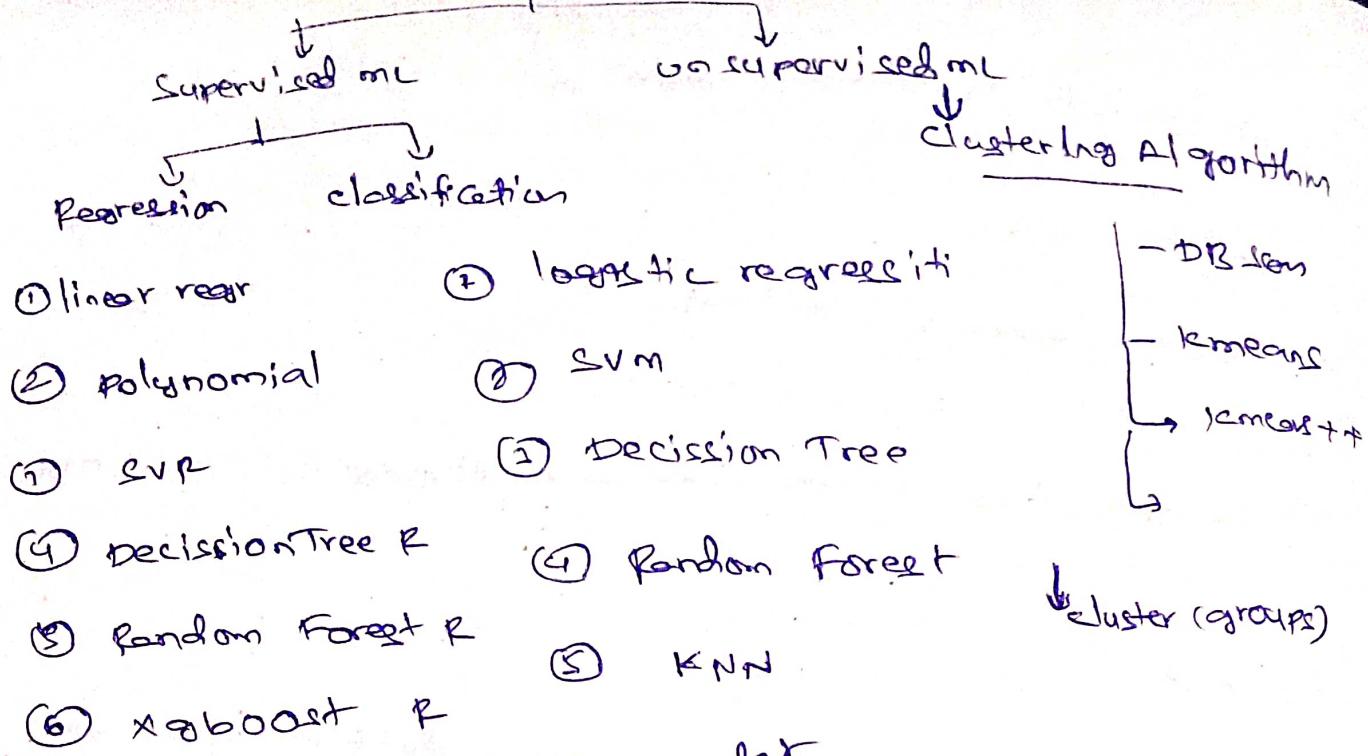
Artificial Intelligence

If is creating an application where it performs all its task without any Human interfaces

Human interfaces

Chatbots → AI Chatbot

Ai → subset is ml → subset is DL
ml provides starts tools to analyse, visualize, perform prediction and other task with the help of data



Supervised ML

Independent feature

Dependent feature

No. of plants hrs	No. of study hrs	Pass/fail
1	1	0
2	2	0
3	5	1

⑩ or ⑪ Classification

Example ↗

Regression

Degree	Exp	Salaries
R.Tech	1	50K
PhD	2	100K

Flight price prediction → Regression

Algerian Fire forest → Classification

Air Quality Index → Regression

Tom (Run/Not) → Classification

Buy Day of the person → Classification

S, M → Day

Unsupervised one

Age	Salary	Spending	Score (1-10)
24	70K	1	
26	100K	9	
21	-	9	
25	120K	2	

(Customer Segmentation)

XX

YY

- * offer 10% or 20%

↑ sales increase 20%

contains
(1 Independent, 1 Dependent) feature

Simple Linear Regression

- * we find best fit line in two variables \rightarrow linear relationship

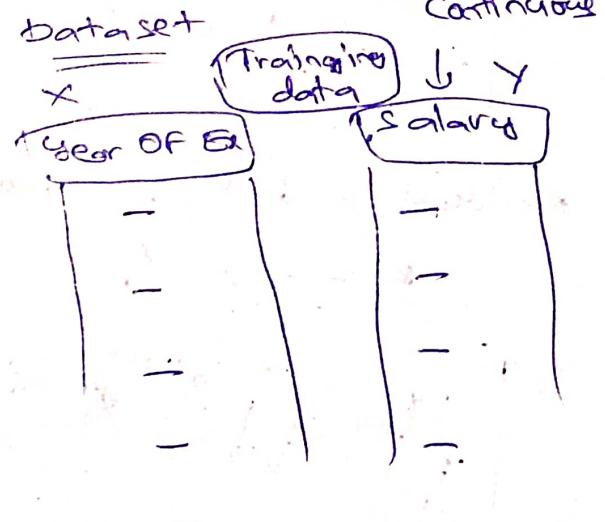
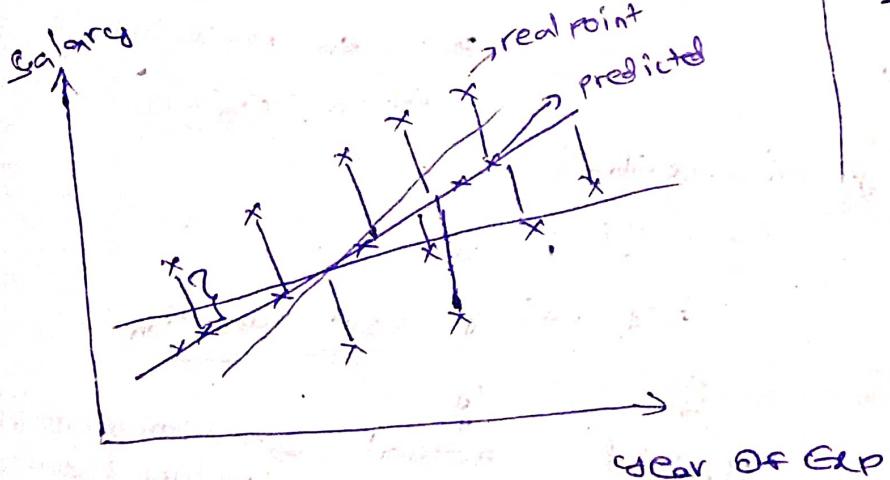
Affine

Aim \rightarrow To Create a model

model \rightarrow Year of Experience & Salary \Rightarrow input

Predict \rightarrow Salary Based on I/o year

Take 1 Independent & 1 dependent



- * minimal main aim find best fit line should be minimal
- * difference b/w actual and predicted (Error, Residuals)
- * minimal distance model is trained well
- * summation should be minimal (less distance)

Equation of straight line

* best fit line is equation of straight line

$$y = \theta_0 + \theta_1 x$$

or

$$y = \beta_0 + \beta_1 x$$

or

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

Intercept

slope

(C)

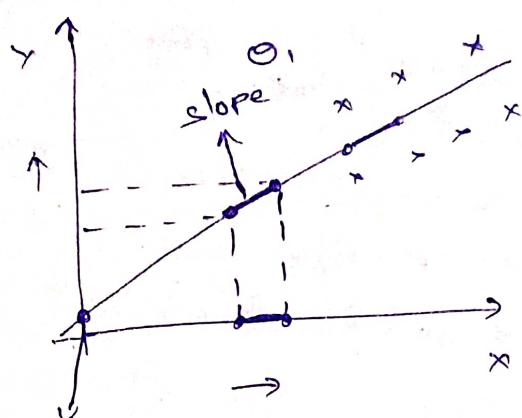
$f(mx)$

unit moment,

x axis is what

unit momenti

y axis is



Intercept θ_0

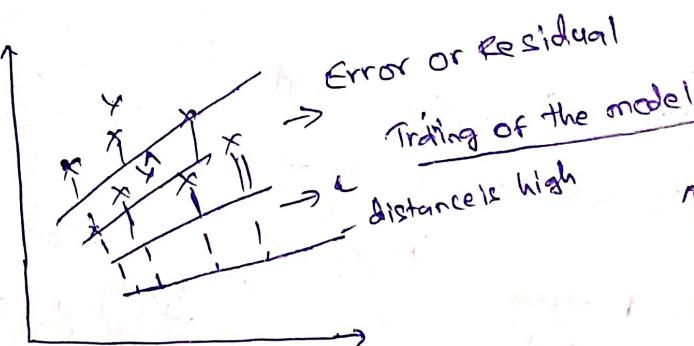
* Intercept and slope change
for best fit line

Example of fitting

Gap
0

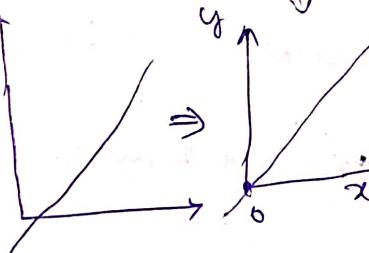
Salary
3.25 lakh

$\theta_0, \theta_1 \rightarrow$ Value



Intercept = 0

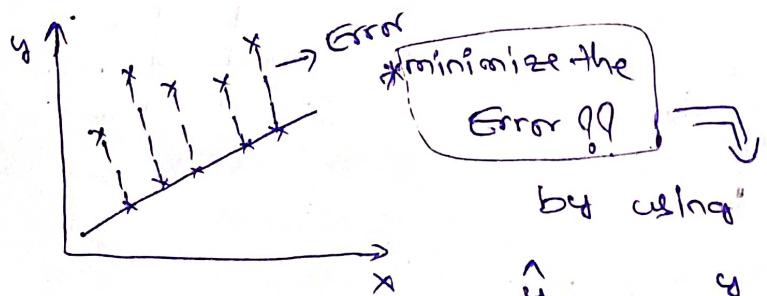
$$\boxed{\theta_0 = 0}$$



* when pass origin
then $\theta_0 = 0$

residual = distance between y and \hat{y}

$$r = y - \hat{y}$$



by using cost function

Cost function :

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

Actual \Rightarrow minimize error /
reduce distance /
decrease
 \Downarrow mean square
error
 \Downarrow \Downarrow

\hat{y}_i = Predicted value
 mse = mean square error
 n = no. of data points
 y_i = Observed values

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

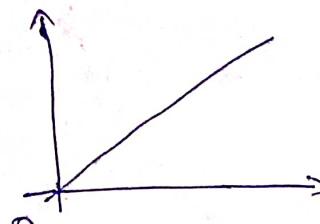
Final Aim:

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

↓ m i=1

change θ_0 and θ_1 to minimize $J(\theta)$

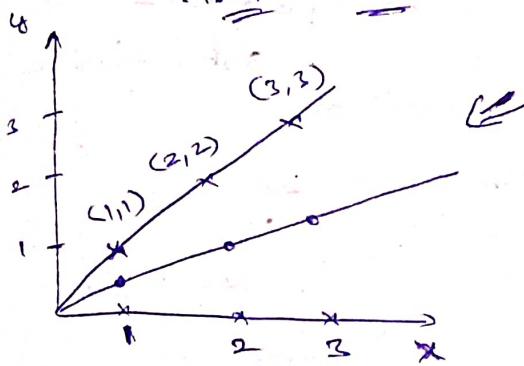
$$\Rightarrow h_\theta(x) = \theta_0 + \theta_1 x$$



lets consider

$$\begin{aligned} \theta_0 &= 0 & h_\theta(x) &= \theta_0 + \theta_1 x \\ \Rightarrow h_\theta(x) &= \theta_1 x \end{aligned}$$

Train Graph



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

$$\theta_1 = 1$$

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

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

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

Train

x	y
1	1
2	2
3	3

$$\theta_1 = 0.5$$

m = total no. of elements value

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

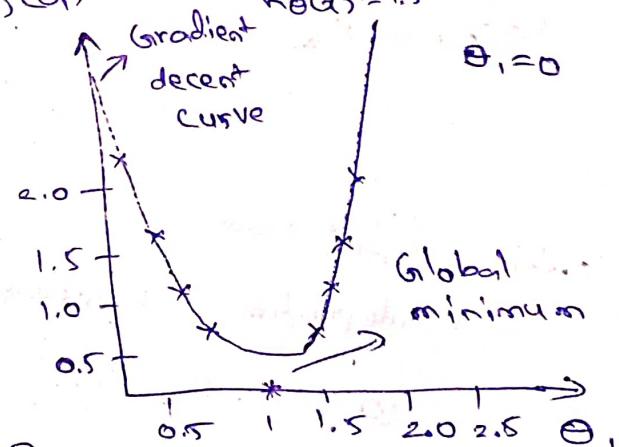
$$h_\theta(x) = 1.0 \quad x=2$$

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

$$\theta_1 = 0$$

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

$$\begin{aligned} &= \frac{1}{3} \left[(1-1)^2 + (2-2)^2 + (3-3)^2 \right] \\ &\Rightarrow \frac{1}{3} (0+0+0) \\ &\Rightarrow 0 \quad \text{it is zero} \\ &\quad \text{so decrease} \downarrow \downarrow \end{aligned}$$



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

$$= \frac{1}{3} [0.25 + 1 + 2.25]$$

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

$$J(\theta_1) = \frac{1}{3} \left[(0-1)^2 + (0-2)^2 + (0-3)^2 \right]$$

$$J(\theta_1) \approx 4.66$$

Convergence Algorithm Optimize the Change of θ , value?

Repeat until Convergence

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

learning rate

↓

θ_0, θ_1

$$\theta_j = \theta_j - \alpha (-ve)$$

$$= \theta_j + \alpha$$



$$\theta_j = \theta_j - \alpha (+ve \text{ slope}) \rightarrow \text{decrease } \theta_j$$

learning $\alpha \Rightarrow$ speed of slope $\alpha = 0.0000001$

$$\alpha = 0.001$$

↓
MAE

↓
RMSE

Tone \Rightarrow V

↳ Cost function

$x_1 \quad x_2 \quad \dots \quad x_n \quad y$

No of rooms	City	Room size	Price
-------------	------	-----------	-------

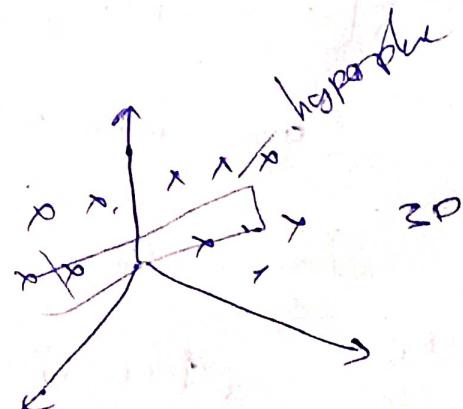
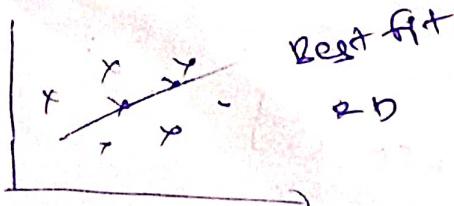
Cohen simple linear

1 independent 1 dependent

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

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

↳ multiple features



hyperplane

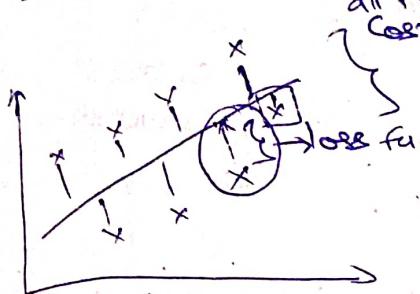
Simple linear regression

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

multiple linear regression

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

Convergence algorithm



all points
Cost function

m = no. of data points

(MSE) is

Cost function: all points applied
every point we calc different
summation

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

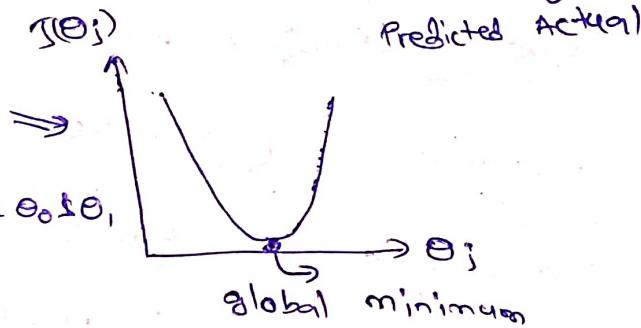
* single point * each point * every point
Loss function := $(h_{\theta}(x^{(i)}) - y^{(i)})^2 \Rightarrow (\hat{y}_{(i)}, y_{(i)})$

$$j = 1 \text{ to } m$$

* repeat until Convergence

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

* update θ_0 & θ_1



slope

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

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

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

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

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

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

Example

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

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

$$\begin{aligned} \frac{\partial}{\partial x} (x+1)^2 &= 2(x+1) \\ &\times (1+0) \\ &\Rightarrow 2(x+1) \end{aligned}$$

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

$(1+0)$

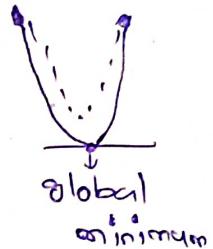
$$\hat{y} = \theta_0 + \theta_1 x$$

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

$$\Rightarrow \frac{\partial}{\partial \theta_1} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})^2 * [x] \quad \begin{array}{l} \text{learning rate} \\ \downarrow \\ \alpha = \text{Speed of convergence} \end{array}$$

Repeat until Convergence:

updated
 θ_0, θ_1



$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})$$

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

Cost function types

① MSE

② MAE

③ RMSE

linear regression
 $(a-b)^2 = a^2 - 2ab + b^2$

① MSE $\frac{1}{n} \sum (y - \hat{y})^2$

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

$$ax^2 + bx + c = 0$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 \quad \begin{array}{l} \text{predicted value} \\ \downarrow \end{array}$$

\downarrow (quadratic equation)

θ_0, θ_1

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

Advantages:

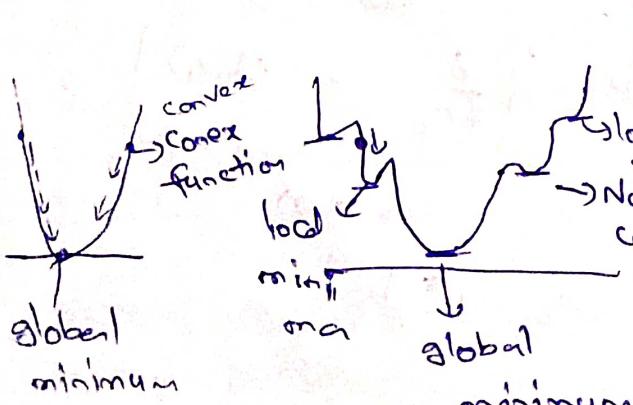
① this equation is differentiable

② this equation also has one global minimum

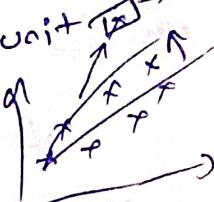
Disadvantages:

① this is not robust to outliers

cost fun \uparrow error.



② Penalizing the error changing unit \rightarrow outlier



remove outlier

Exp \rightarrow dependent feature
Salaries $(y - \hat{y})^2$
 \downarrow Salary (Lakh INR) $(\text{Lakh})^2$

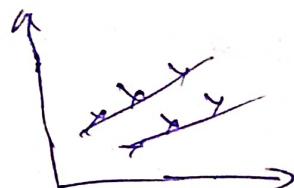
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

\Rightarrow Cost function \downarrow
 reduce cost function

② Mean Absolute Error (MAE)

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

absolute
 $|2-5| = 1-2| = \frac{3}{2}$



Advantage

① Robust to Outliers

② It will also be in the same unit

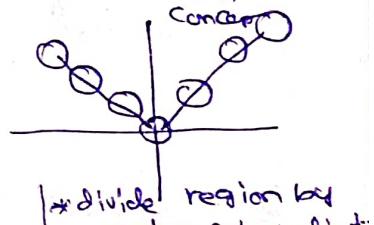
Disadvantage

③ Convergence usually

takes more time \Rightarrow optimization is a complex task

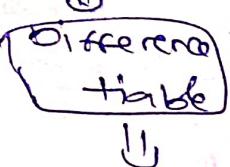
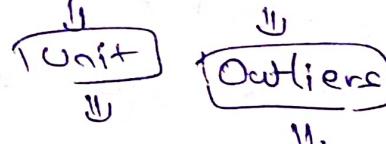


sub gradient descent



divide region by subgradient

* Time Consuming



④ Huber loss

⑤ RMSE

$$\sqrt{\text{MSE}}$$

* Huber loss you want to use the huber loss any time you feel that you need to balance between giving outliers some weight, but not too much for case where outliers very important to you use mse you don't care about outliers.

RMS E

* RMSE has the benefit of penalizing large errors more so can be more appropriate in some cases.

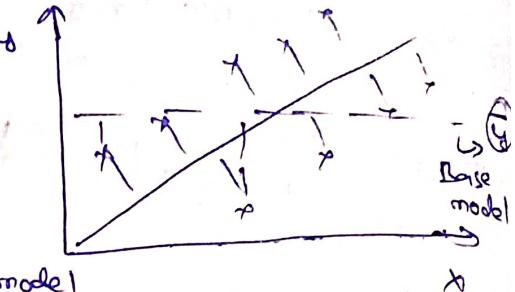
Performance metrics :-

① R Squared :-

① R Squared :-

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

$$= 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \rightarrow \begin{cases} \text{low value if model is fit well} \\ \text{high value if model is not fit well} \end{cases}$$



\bar{y} = Average of y

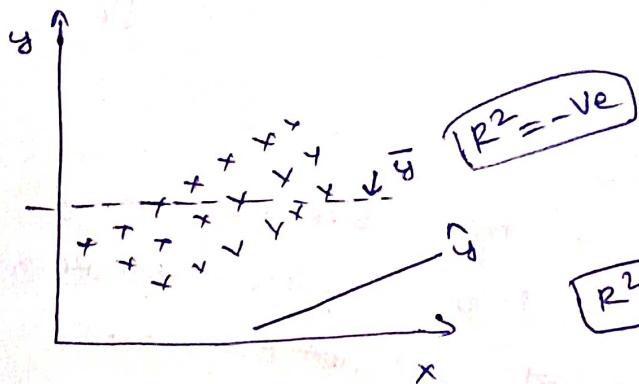
$$R^2_{\text{Squared}} = 1 - \frac{\text{small number}}{\text{Bigger number}} \rightarrow \begin{cases} R^2 \leq 1 \\ \Rightarrow 0.85 \rightarrow 85\% \text{ Accuracy} \end{cases}$$

$\Rightarrow 0.75\%$

$\Rightarrow 75\% \text{ Accuracy}$

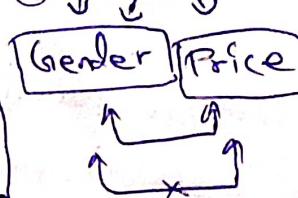
$\Rightarrow R^2_{\text{Squared}} \text{ is -ve}$
measured model is very bad

Performance of the model that you have created



- * add feature
- * add city location R^2 increase again
- * add no. of rooms R^2 increase again

④ ↘ ↗



① R^2 65%

$Adj R^2 = 63\% P = ?$
 $Adj R^2 = 73\% P = ?$

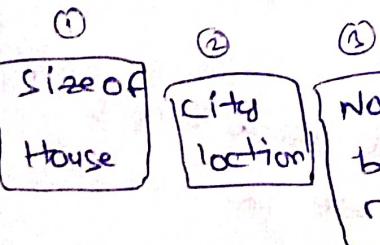
② $R^2 \geq 75\%$

③ $R^2 88\%$

④ $R^2 90\%$

$Adj R^2 85\%$

② Adjusted R Squared :-



* R^2 is increased in slightly

* $Adj R^2$ is more than R^2

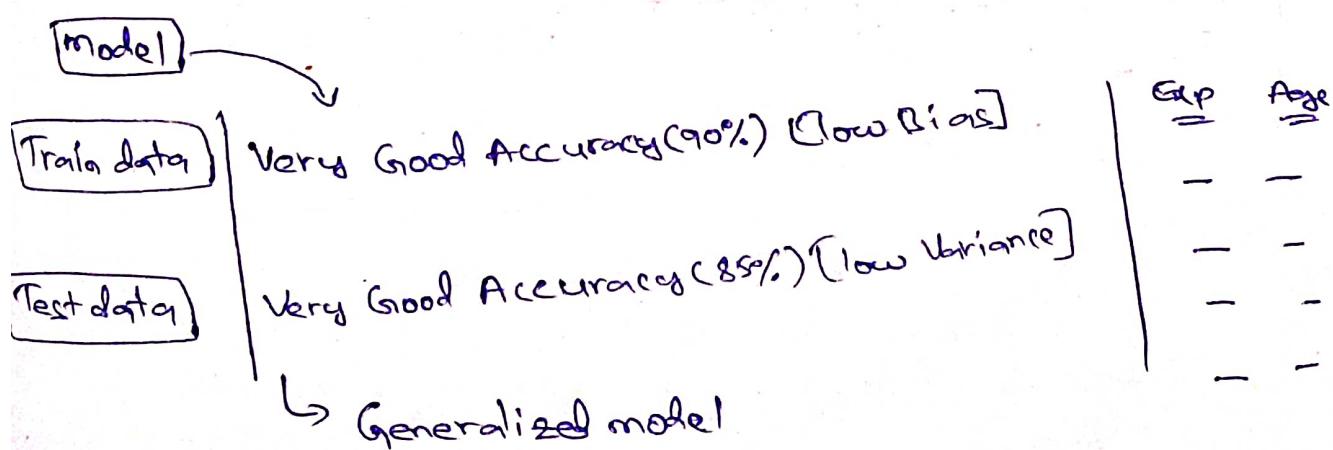
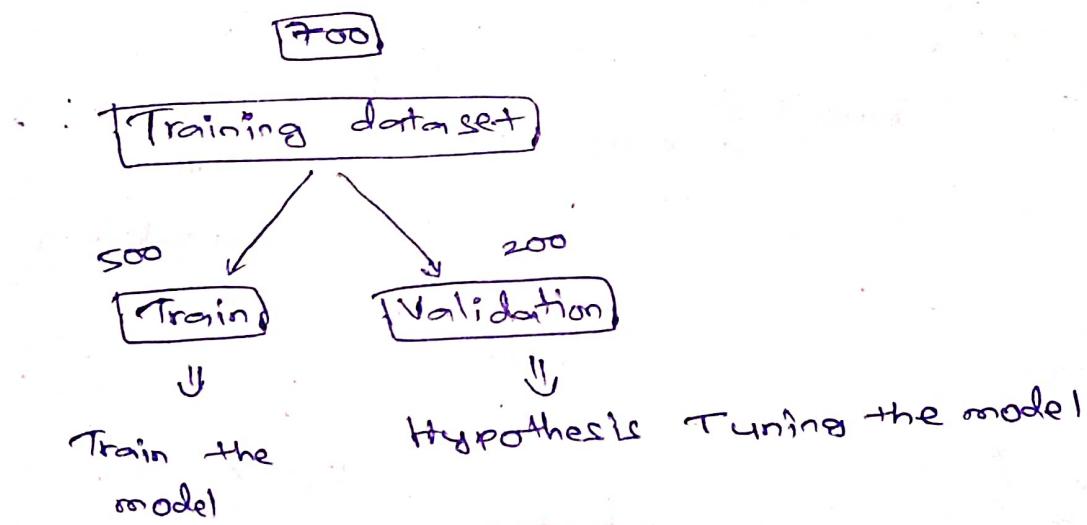
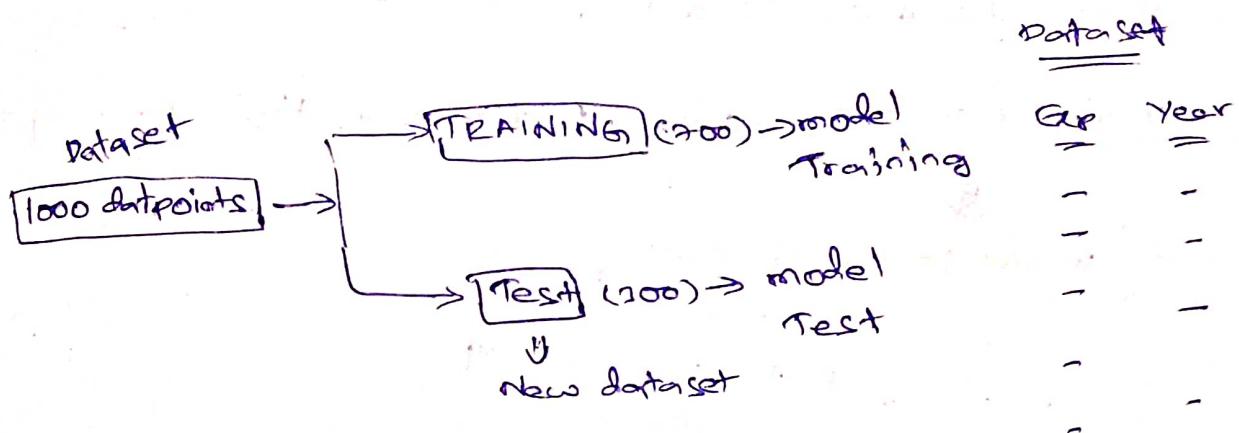
$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(n-1)}{N-p-1}$$

n = No. of data points

p = No. of independent feature

* add unnecessary column
in R^2 incre but
adj R^2 is decreased
(adj R^2 best)

Overfitting And Underfitting (Bias And Variance)



Train Very Good Accuracy (90%) [low Bias]

Test Bad Accuracy (50%) [High Variance]



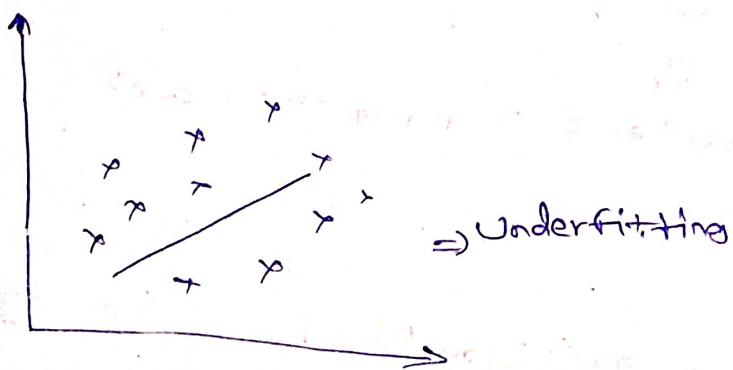
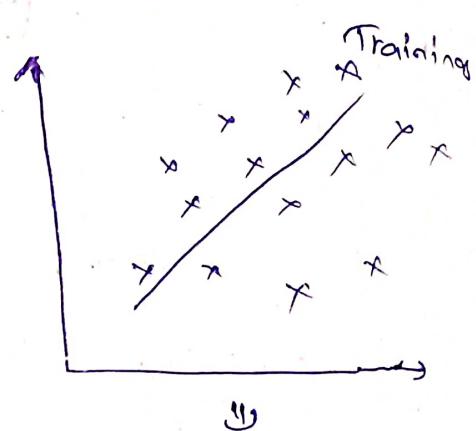
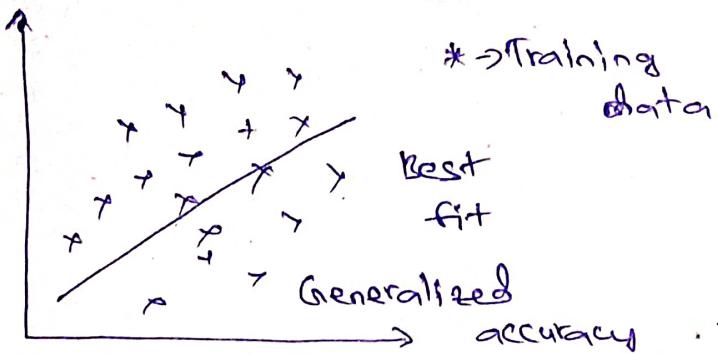
Overfitting

model accuracy is low [High Bias]

model accuracy is low [Low or High Variance]

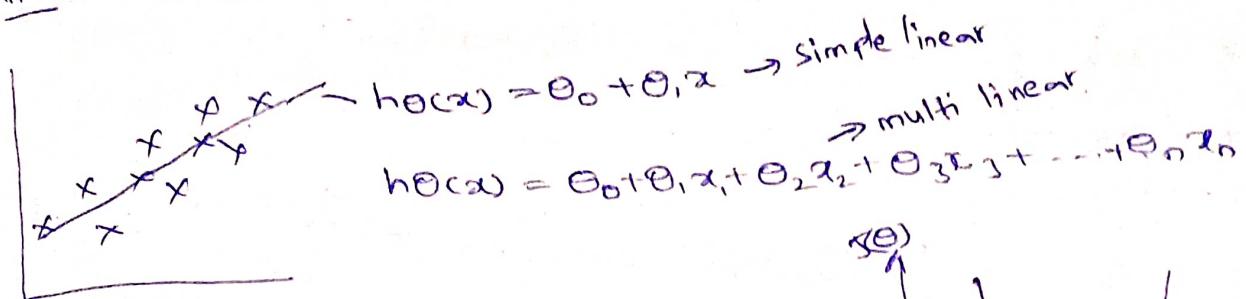


model is underfitting



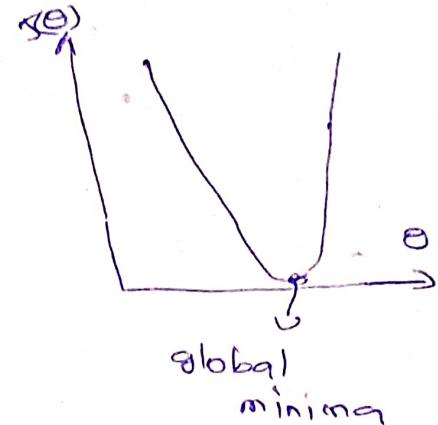
Ridge Regression , Lasso Regression , Elasticnet regression

Linear regression



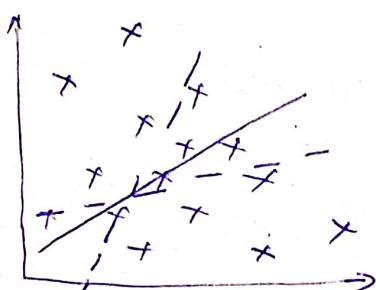
cost function = $\frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$

mean square error



① Ridge Regression \Leftrightarrow (L₂ regularization)

* used to reduce Overfitting * cost function = 0



Overfitting

(local bias)
 Training data \rightarrow Accuracy $\uparrow\uparrow = 95\%$.

Test data \rightarrow Accuracy $\downarrow\downarrow = 60\%$
 (high variance)

hyper parameter



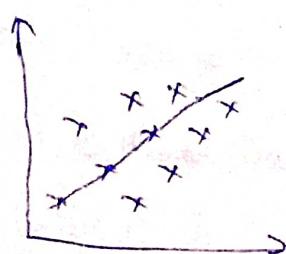
Cost fun. $\frac{\partial}{\partial \theta} \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{i=1}^n (\text{slope})^2$

$\hookrightarrow [\theta_0^2 + \theta_1^2 + \theta_2^2 + \dots + \theta_n^2]$

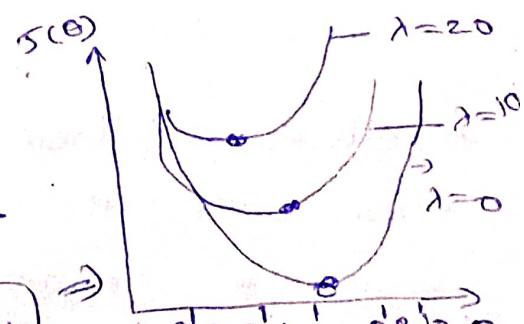
$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$= \theta_0 + 0.95x_1 + 0.82x_2 + 1.5x_3$$

$$= \theta_0 + 0.25x_1 + 0.70x_2 + 0.77x_3$$



$\lambda \uparrow \quad \theta \downarrow$



cost fn = $0 + \lambda (\text{slope})^2$
 $= \# \text{ vers } \downarrow \downarrow$

② Lasso regression is [L_1 , norm]
 L_1 regularization

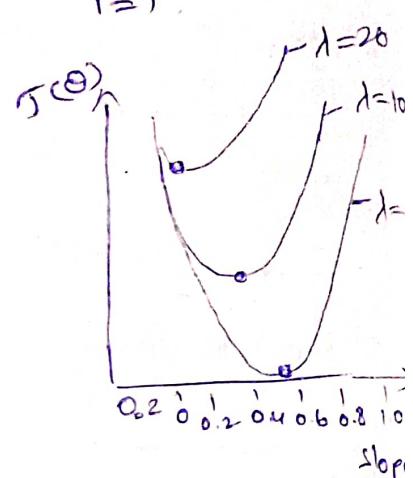
- * used to reduce the features
- * feature selection

$$\text{Cost fn: } \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{i=1}^m |\text{slope}|$$

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

$$= \theta_0 + \theta_1, 0.54x_1 + 0.27x_2 + [0.10x_3]$$

↓
removed
feature



③ Elastic [L1 and L2 norm] :

$$\text{Cost fn} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda_1 \sum_{i=1}^m (\text{slope})^2 + \lambda_2 \sum_{i=1}^m |\text{slope}|$$

↑
ridge
↑
lasso

- * Elasticnet linear regression uses the penalties from both the lasso and ridge techniques to regularize regression models.
- * both methods by learning from their shortcomings to improve the regularization of statistical models