

# Machine Learning

## 7-days Machine Learning Live Sessions

### Agenda :-

- ① Introduction to ML (AI vs ML vs DL vs DS)
- ② Supervised ML vs Unsupervised ML
- ③ Linear Regression (Math's and Geometric intuition)
- ④  $R^2$  and Adjusted  $R^2$
- ⑤ Ridge and Lasso Regression.

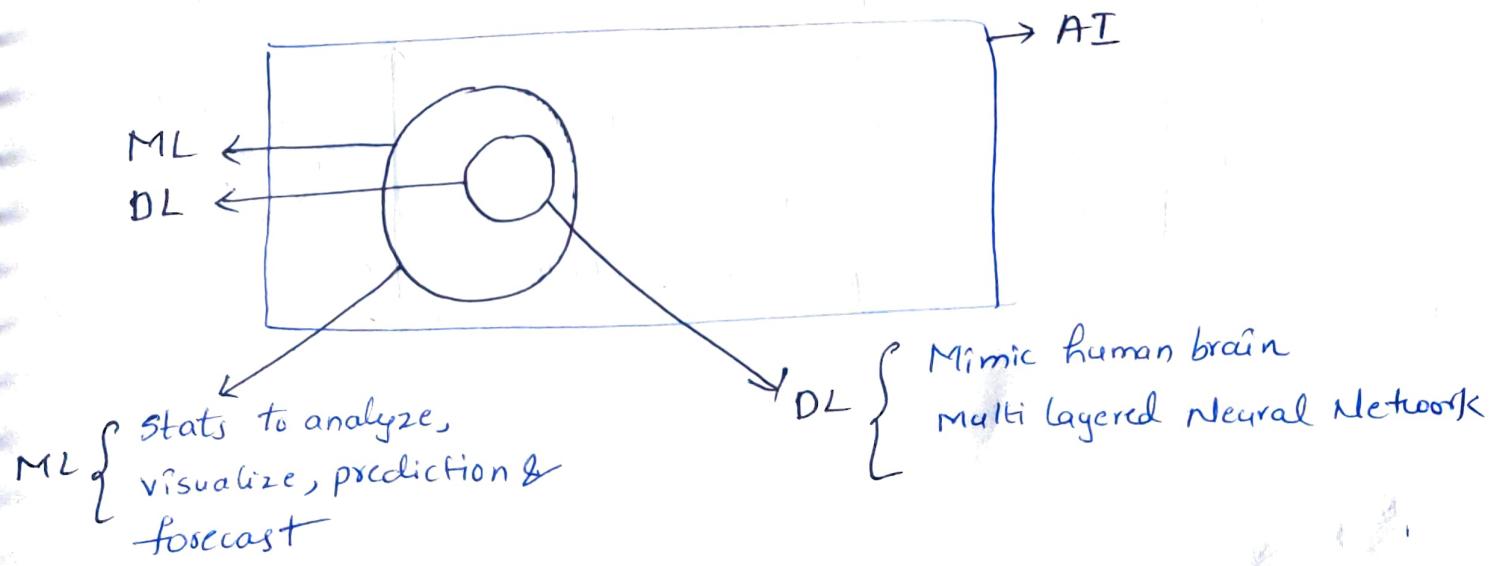
### ① AI vs ML vs DL vs DS

→ AI application is able to do its own task without any human intervention

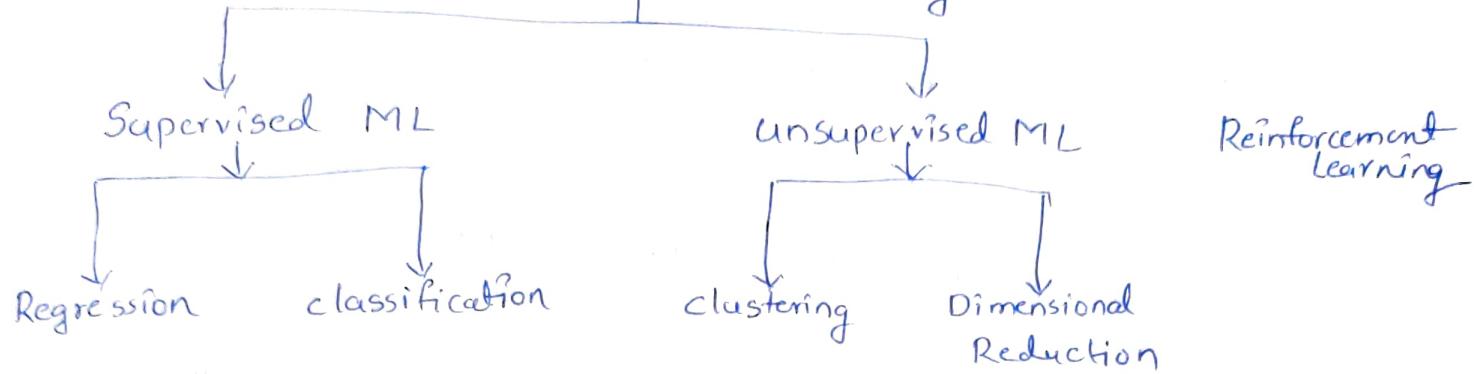
Eg t Netflix App → Recommendation

Action  
Comedy  
Romance

Amazon.in → Recommendation



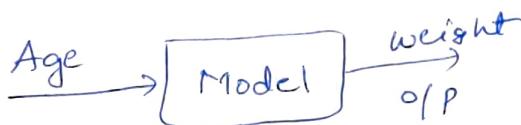
# Machine & Deep Learning



Supervised ML :-

- \* Labeled data
- \* Independent and dependent features

↑ I/P Age	↓ O/P weight
24	62
25	63
21	72
27	62



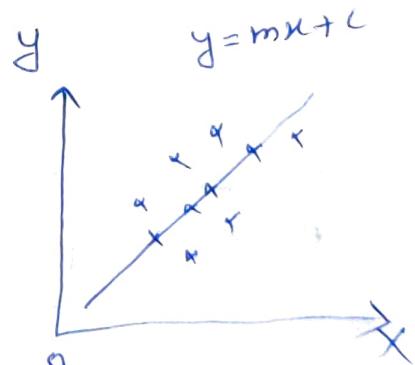
Independent feature : Age

Dependent feature : weight

(i) Regression problem :-

Age	weight	→ O/P
24	62	
25	63	
21	71	
27	62	

↓  
Continuous variable  
↓  
Regression problem



## (ii) classification problem

\* output variable  $y$  Discrete (categorical)

No. of hours	No. of play hours	No. of house Sleep	O/P
-	-	-	P
-	-	-	F
-	-	-	F
-	-	-	P

only 2 outputs  $\Rightarrow$  Binary classification

More than 2 outputs  $\Rightarrow$  Multiclass classification.

Unsupervised ML is

\* No output or dependent variable

## (i) clustering

Salary	Age
20k	20
30k	25
40k	35
50k	42

clustering  $\rightarrow$  customer segmentation



## (ii) Dimensionality Reduction

1000 features  $\xrightarrow{\text{Convert to}}$  lower dimensions

Eg (100 features)

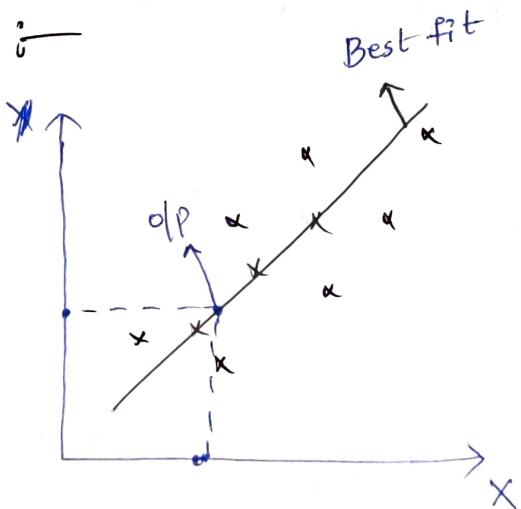
## Supervised ML

- Linear Regression
- Ridge & Lasso
- Logistic Regression.
- Decision Tree
- AdaBoost
- Random Forest
- Gradient Boosting
- XGBoost
- Naive Baye's
- SVM
- KNN

## unsupervised ML

- k-Means
- DBSCAN
- k Nearest Neighbour cluster
- PCA
- LDA

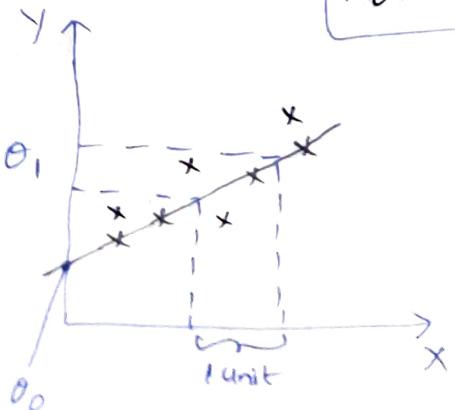
## Linear Regression



y is a linear function of x

Suppose equation of a straight line

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



$\theta_0$  = Intercept { where best line meet at y-axis }

$\theta_1$  = scope or coefficient

for 1 unit in x-axis vary in y-axis

Best fit line  $\hat{y}$  — The distance between actual values to line is minimum.

Cost function  $J$  —

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

→ cost function

called as squared error function

$h_\theta(x^{(i)})$  — Predicted

$y^{(i)}$  — Actual.

$\therefore \frac{1}{2m}$  used to derivation purpose.

what we need to solve

Minimize

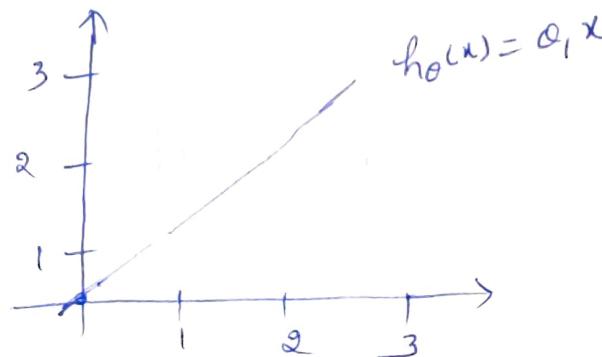
$$J(\theta_0, \theta_1)$$

$$\theta_0, \theta_1$$

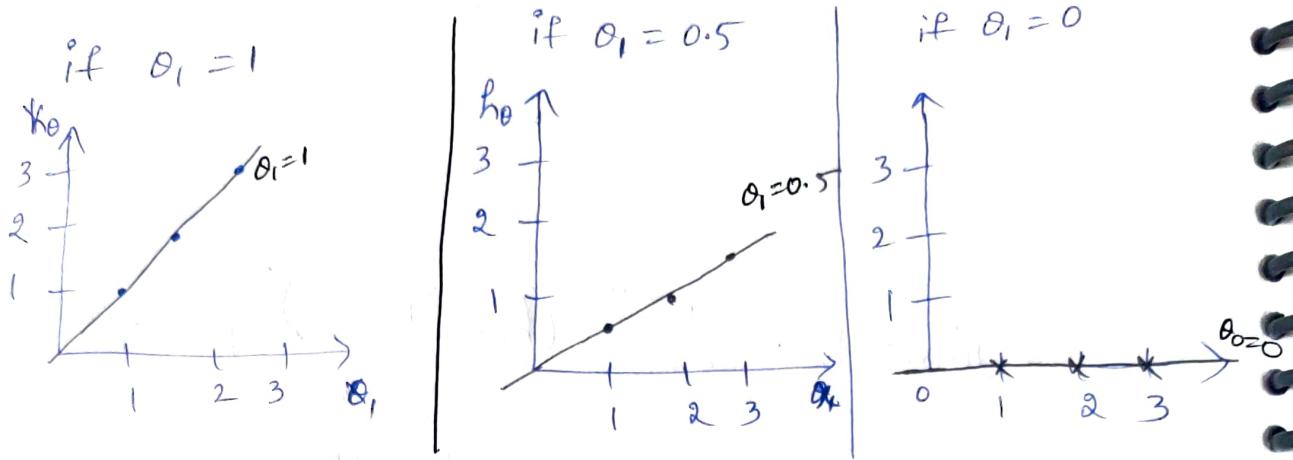
$$= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

$$\text{if } \theta_0 = 0 \implies h_\theta(x) = \theta_1 x$$



Suppose data points  $(1,1) (2,2) (3,3)$



if  $\theta_1 = 1$

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

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

$$= \frac{0+0+0}{2m}$$

$$= 0$$

$\text{if } \theta_1 = 1 \Rightarrow J(\theta_1) = 0$

if  $\theta_1 = 0.5$

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

$$= \frac{1}{2 \times 3} [0.25 + 1 + 2.25]$$

$$= 0.58$$

$\text{if } \theta_1 = 0.5 \Rightarrow J(\theta_1) = 0.58$

If  $\theta_1 = 0$

$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2 \times 3} [(0-1)^2 + (0-2)^2 + (0-3)^2] \\ &= \frac{1}{6} [1+4+9] \\ &= 2.5 \end{aligned}$$

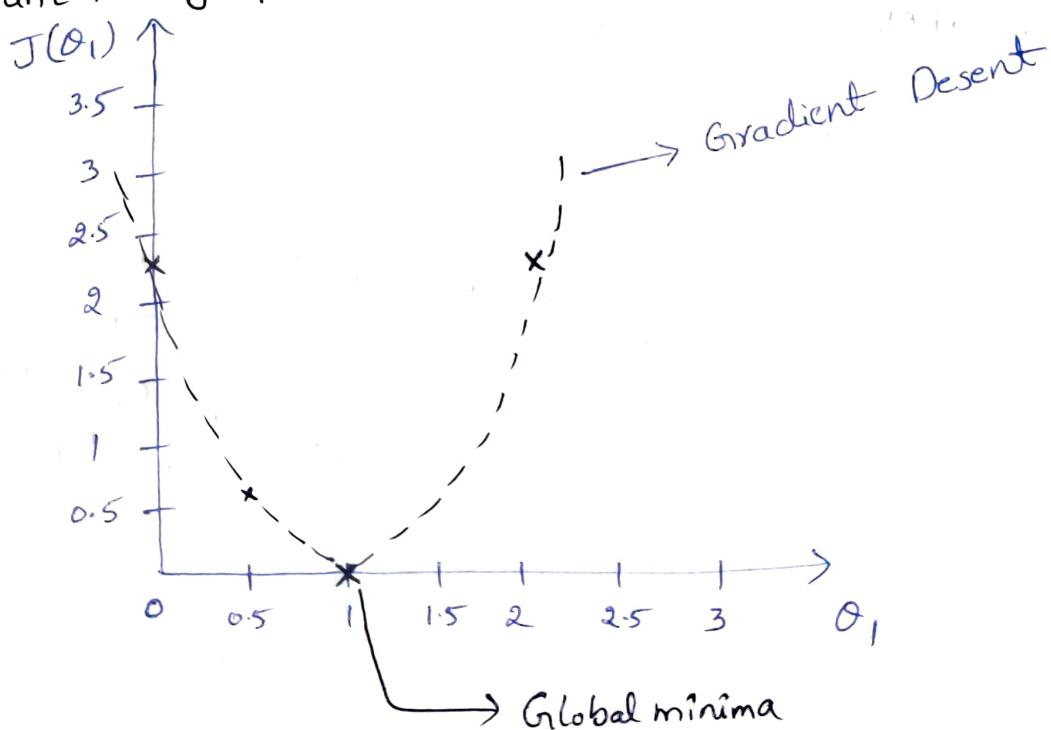
$$\boxed{\text{if } \theta_1 = 0 \Rightarrow J(\theta_1) = 2.5}$$

If  $\theta_1 = 2$

$$\begin{aligned} J(\theta_1) &= \frac{1}{2 \times 3} [(2-1)^2 + (4-2)^2 + (6-3)^2] \\ &= \frac{1}{6} [1+4+9] \\ &= 2.5 \end{aligned}$$

$$\boxed{\text{if } \theta_1 = 2 \Rightarrow J(\theta_1) = 2.5}$$

Cost function graph

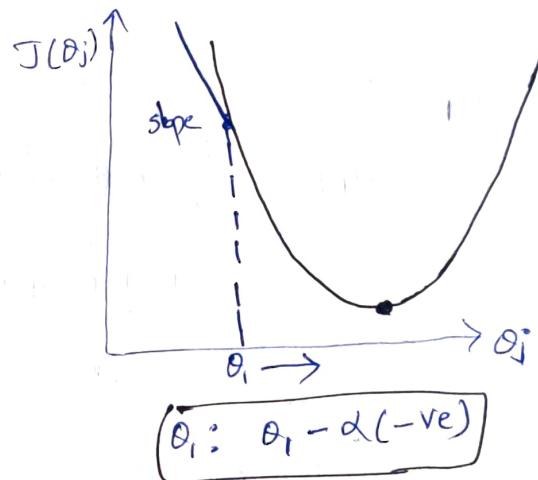
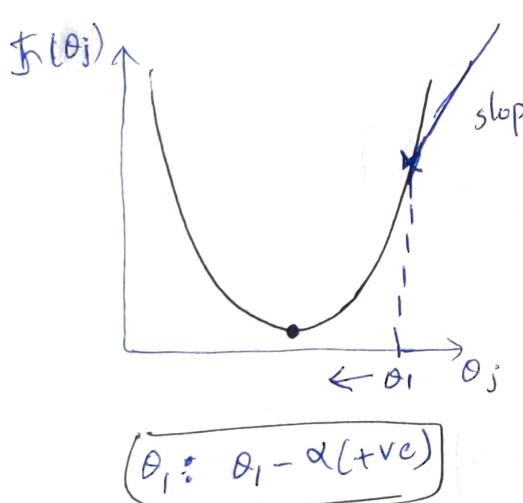


## Convergence algorithm

Repeat until Convergence

$$\left. \begin{array}{l} \\ \end{array} \right\} \quad \theta_j = \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$

$\alpha$  = Learning Rate  $\Rightarrow$  how speed the value  $\theta_j$  changes



if graph likes



## GRADIENT DESCNT Algorithm

$$\left\{ \theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} (J(\theta_0, \theta_1)) \right.$$

}

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

with r.t to  $\theta_0$

$$\begin{aligned} \frac{\partial}{\partial \theta_0} (J(\theta_0, \theta_1)) &= \cancel{\frac{1}{2m}} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \end{aligned}$$

w.r.t  $\theta_1$

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

Repeat until convergence

{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

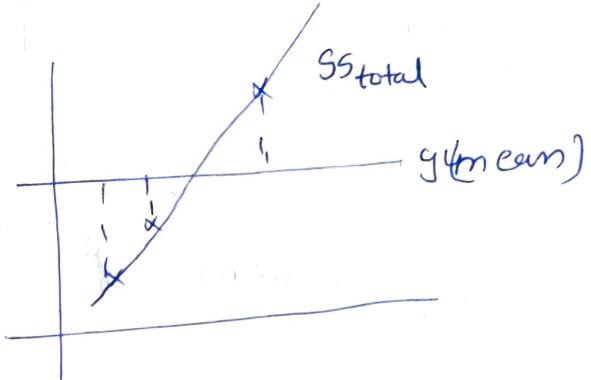
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

}

$R^2$  and Adjusted  $R^2$  :-

\*  $R^2$  is performance metrics.

$$R^2 = 1 - \frac{SS_{\text{Res}}}{SS_{\text{total}}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$



if non-correlated features added, then  $R^2$  increased. This should be avoided., so using Adjusted  $R^2$

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

N = No. of samples

P = Features or predictors

Eg:-

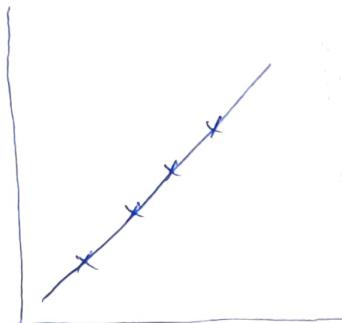
$$\left\{ \begin{array}{l} P=2 \Rightarrow R^2 = 90\% \quad \text{Adjusted } R^2 = 86\% \\ P=3 \Rightarrow R^2 = 91\% \quad \text{Adjusted } R^2 = 82\% \end{array} \right.$$

Day-2

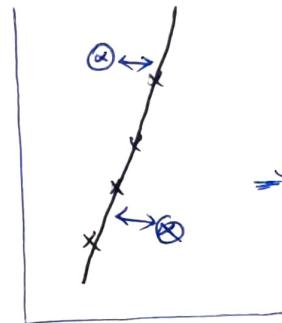
- ① Ridge & Lasso Regression
- ② Assumption of Linear Regression
- ③ Logistic Regression
- ④ Confusion matrix
- ⑤ practical implementation.

Ridge & Lasso Regression :-

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



Training data  
fitted very well



test data  
not able generalized

⇒ This is  
overfitting

overfitting :-

Model performs well on - Training data

fail to perform well on - Test data

Low Bias

High Variance

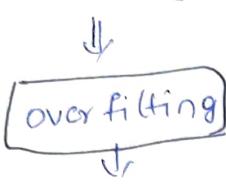
underfitting :-

Model accuracy bad with training data ⇒ High Bias

Model accuracy also bad with test data ⇒ High Variance

### Model 1

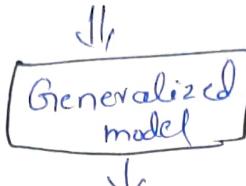
Train Accuracy = 90%  
Test Accuracy = 80%



Low Bias, High Variance

### Model 2

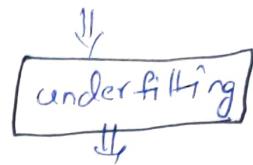
Train Accuracy = 92%  
Test Accuracy = 91%



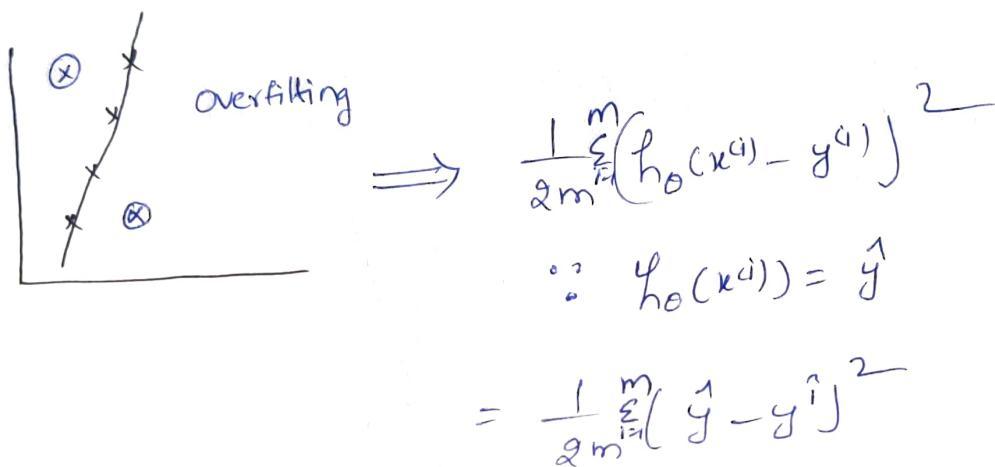
Low Bias,  
Low Variance

### Model 3

Train Accuracy = 70%  
Test Accuracy = 65%



High Bias,  
High Variance



Now introduce Ridge by adding  $\lambda(\text{slope})^2$  to Cost function.

Ridge [L2 Regularization] :-

$\lambda$  - hyperparameter  
iteration  $\rightarrow$  has many times  $\theta$  value need to change

$$\text{Cost function} + \lambda(\text{slope})^2$$

$$= \frac{1}{2m} \sum_{i=1}^m (\hat{y} - y^{(i)})^2 + \lambda(\text{slope})^2$$

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

- ⊗ By using L2 regularization, we can prevent overfitting and generate generalized model.

## Lasso [ L1 Regularization ]

$|\text{slope}| \Rightarrow$  helps to feature selection

Cost function +  $\lambda |\text{slope}|$

$$= \frac{1}{2m} \sum_{i=1}^m (\hat{y} - y^i)^2 + \lambda |\text{slope}|$$

Suppose

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

$$|\text{slope}| = |\theta_0 + \theta_1 + \theta_2 + \theta_3 + \dots + \theta_n|$$

(i) preventing overfitting

(ii) helps to feature selection (Remove unwanted features)

## Ridge Regression (L2 Regularization)

$$\text{Cost function} = (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda (\text{slope})^2$$

Purpose : Preventing overfitting

## Lasso Regression (L1 Regularization)

$$\text{Cost function} = (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda |\text{slope}|$$

Purpose : (i) prevent overfitting  
(ii) feature selection.

## Assumptions of Linear Regression

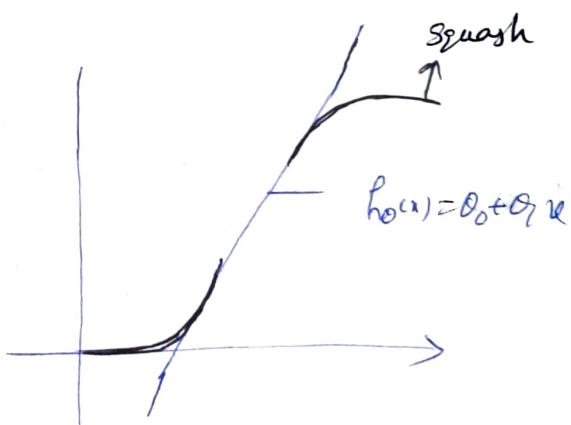
- (i) Gaussian/Normal Distribution  $\rightarrow$  Model will get trained well
- (ii) Standardization of scaling data  $\Rightarrow$  z-score  $[\mu=0, \sigma=1]$
- (iii) Linearity
- (iv) Multi Colinearity.

## Logistic Regression:— [classification]

Decision boundary      Logistic Regression

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

$$h_{\theta}(x) = \theta^T x$$



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

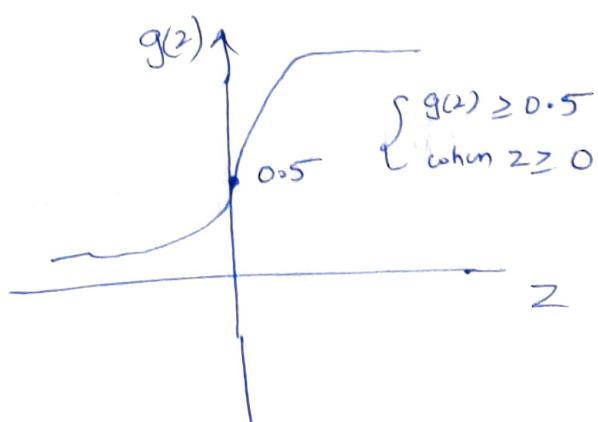
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x)$$

$$h_{\theta}(x) = g(z)$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-z}} \rightarrow \begin{array}{l} \text{Sigmoid} \\ \text{or} \\ \text{logistic function} \end{array}$$

$\Downarrow$

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1)x}}$$



Training set :-

$$\{(x^1, y^1), (x^2, y^2), (x^3, y^3), \dots, (x^n, y^n)\}$$

$y \in \{0, 1\} \rightarrow 20/\text{ps} \Rightarrow \text{Binary classification.}$

$$h_\theta(z) = \frac{1}{1+e^{-z}}$$

$$\therefore z = \theta_0 + \theta_1 x$$

if assume

$$\theta_0 = 0$$

↓,

$$z = \theta_1 x$$

④ change parameter  $\theta_1$  to get best fit line.

Cost function

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

logistic Regression = Replace  $h_\theta(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}}$  → this cannot be used as cost function  
→ This has local minima problem.

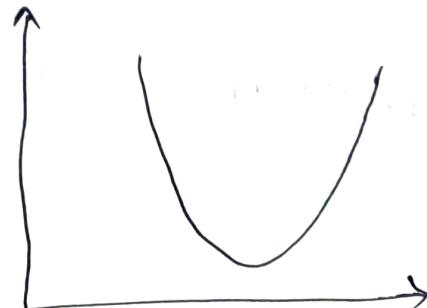
$\frac{1}{1+e^{-z}}$  is a non-convex function.

→ lot of local minima, difficult to get global minima

Non-Convex



Convex

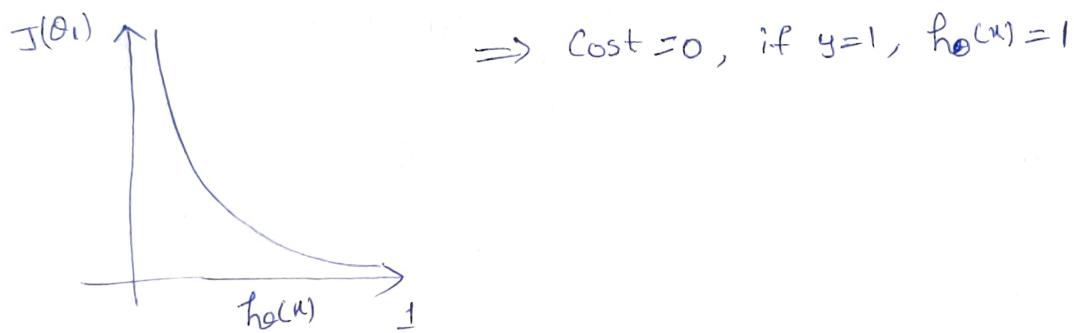


In order to solve local minima problem, the logistic cost function is

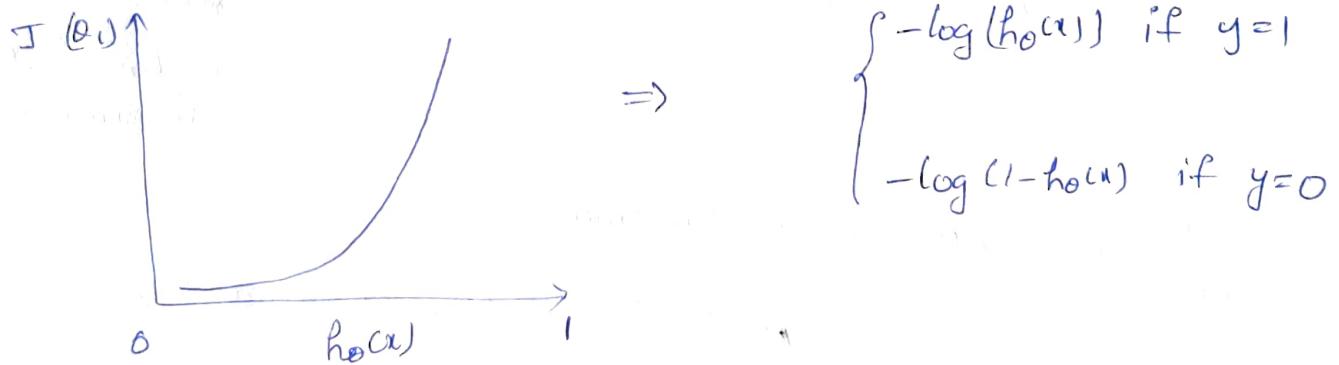
$$J(\theta_0) = \begin{cases} -\log(h_{\theta}(x)) & y=1 \\ -\log(1-h_{\theta}(x)) & y=0 \end{cases}$$

$$\therefore h_{\theta}(x) = \frac{1}{1+e^{-\theta_0(x)}}$$

if  $y=1$



if  $y=0$



$$\boxed{\text{Cost}(h_{\theta}(x^{(i)}), y) = -y \log(h_{\theta}(x^{(i)})) - (1-y) \log(1-h_{\theta}(x))}$$

Substitute  $y=0, y=1$ , you will get before one.

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

$$h_\theta(x^i) = \frac{1}{1+e^{-\theta_0 x^i}}$$

Repeat until Convergence

$$\left. \begin{array}{l} \\ \theta_j := \theta_j - \alpha \frac{\partial J(\theta_0)}{\partial \theta_j} \end{array} \right\}$$

Performance Matrix { classification problem }

		Actual	
		1	0
Predicted	1	TP	FP
	0	FN	TN

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+FN+TN}$$

$$\text{Precision} = \frac{TP}{TP+FP} \quad \{ \text{spam} \}$$

$$\text{Recall} = \frac{TP}{TP+FN} \quad \{ \text{cancer} \}$$

$$F\text{-Score} = (1+\beta^2) \left( \frac{P * Re}{P + Re} \right)$$

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$\text{Recall or TPR or Sensitivity} = \frac{TP}{TP + FP}$$

$$\text{Precision} = \frac{TP}{TP + FN}$$

$$F\text{-Score} = (1+\beta^2) \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{if } \beta = 1 \Rightarrow \frac{2xy}{x+y} \Rightarrow F1 \text{ score}$$

$$\text{if } \beta = 0.5 \Rightarrow (1+0.5^2) \left( \frac{P \times R}{P+R} \right), \quad FP \gg FN$$

$$\text{if } \beta = 2 \Rightarrow FN \gg FP \Rightarrow F2 \text{ score.}$$

Day-3G

① practical Implementation

② Naïve Baye's

③ KNN

## Naïve Baye's :-

- \* used for classification problem
- \* Baye's theorem

$$P(A \text{ and } B) = P(B \text{ and } A)$$

$$P(A) * P(B/A) = P(B) * P(A/B)$$

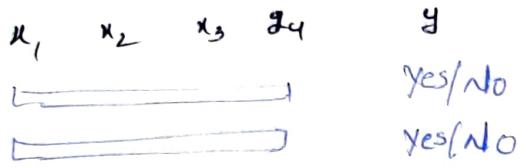
$$P(B/A) = \frac{P(B) * P(A|B)}{P(A)} \implies \text{Baye's theorem}$$

Suppose

$$\underbrace{x_1, x_2, x_3, x_4, \dots, x_n}_{\text{?/P}} \quad \underbrace{y}_{{}^0/\text{P}}$$

$$\begin{aligned} P(y/x_1, x_2, \dots, x_n) &= \frac{P(y) * P(x_1, x_2, x_3, \dots, x_n/y_i)}{P(x_1, x_2, x_3, \dots, x_n)} \\ &= \frac{P(y) P(x_1/y_1) P(x_2/y_2) P(x_3/y_3) \dots P(x_n/y_n)}{P(x_1) P(x_2) P(x_3) \dots P(x_n)} \end{aligned}$$

Assume



$$p(y=yes/x_i) = \frac{p(yes) * p(x_1=yes) * p(x_2=yes) * p(x_3=yes) * p(x_4=yes)}{p(x_1) p(x_2) p(x_3) p(x_4)}$$

$$p(y=No/x_i) = \frac{p(No) * p(x_1=No) * p(x_2=No) * p(x_3=No) * p(x_4=No)}{p(x_1) p(x_2) p(x_3) p(x_4)}$$

Here  $p(x_1) p(x_2) p(x_3) p(x_4) \rightarrow$  is fixed constant, so ignore.

If  $p(y=yes/x_i) = 0.13$  and  $p(y=No/x_i) = 0.05$

$$\geq 0.5 \Rightarrow 1$$

$$< 0.5 \Rightarrow 0$$

Now do Normalization

$$p(yes/x_i) = \frac{0.13}{0.13 + 0.05} = 0.72 = 72\% \quad \left. \right\} 72\%. \text{ Yes best}$$

$$p(No/x_i) = 1 - 0.72 = 28\%$$

## (Q) Dataset

Day	outlook	Temparature	Humidity	wind	Tennis
D1	Sunny	Hot	High	weak	No
D2	Sunny	Hot	High	strong	No
D3	overcast	Hot	High	weak	Yes
D4	Rain	Mild	High	weak	Yes
D5	Rain	Cool	Normal	weak	Yes
D6	Rain	Cool	Normal	Stronge	No
D7	overcast	Cool	Normal	Stronge	Yes
D8	Sunny	Mild	High	weak	No
D9	Sunny	Cool	Normal	weak	Yes
D10	Rain	Mild	Normal	weak	yes
D11	Sunny	Mild	Normal	strong	yes
D12	overcast	Mild	High	Strong	yes
D13	overcast	Hot	Normal	weak	yes
D14	Rain	Mild	High	Strong	No.

Sold

outlook

	Yes	No	p(y)	p(N)
Sunny	2	3	2/9	3/5
overcast	4	0	4/9	4/5
Rain	3	2	3/9	2/5
Total	<u>9</u>	<u>5</u>		

$$P(\text{yes}) = \frac{9}{14}$$

$$P(\text{No}) = \frac{5}{14}$$

Temparature

	Yes	No	P(Y)	P(N)
Hot	2	2	2/9	2/5
Mild	4	2	4/9	4/5
Cool	3	1	3/9	1/5
	<u>9</u>	<u>5</u>		

$$P(\text{yes}) = \frac{9}{14}$$

$$P(\text{No}) = \frac{5}{14}$$

$$P(\text{yes} | (\text{sunny}, \text{Hot})) = \frac{P(\text{yes}) * P(\text{sunny/yes}) * P(\text{Hot/yes})}{P(\text{sunny}) P(\text{Hot})} \quad \text{Ignore}$$

$$= \frac{9/14 * 2/9 * 2/9}{7/14 * 2/9}$$

$$= \frac{3}{5} * \frac{4}{14 * 9}$$

$$= \frac{2}{7 * 9}$$

$$= 2/63 = 0.031$$

$$P(\text{No} | (\text{sunny}, \text{Hot})) = \frac{P(\text{No}) * P(\text{sunny/No}) * P(\text{Hot/No})}{P(\text{Hot}) P(\text{sunny})} \quad \text{Ignore}$$

$$= \frac{5}{14} * \frac{3}{8} * \frac{2}{5}$$

$$= \frac{3}{7 * 5} = \frac{3}{35} = 0.085$$

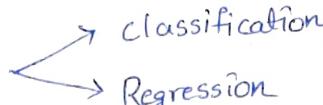
No, do Normalization

$$P(\text{No}) = \frac{0.085}{0.085 + 0.031} = 73\%$$

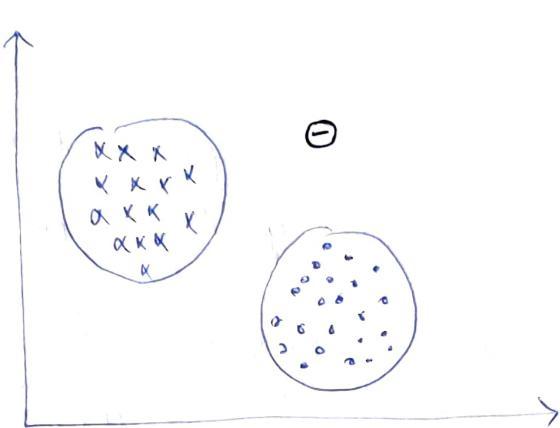
$$P(\text{yes}) = 1 - 73\% = 27\%$$

73%  $\Rightarrow$  Hence answer is No

KNN :-

used to solve  classification  
Regression

## ① classification



Assume = 5  
 $K = 5$

More number of points near to which group that is answer.

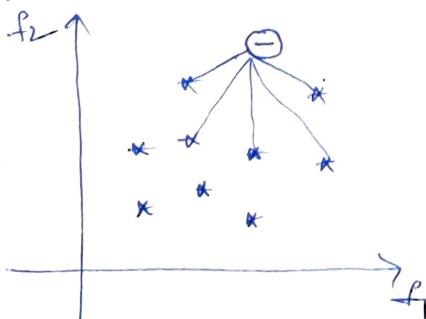
Eucledian distance

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Manhat distance

$$\Rightarrow \{ |x_2 - x_1| + |y_2 - y_1| \}$$

Regression :-



Take the avg of points and that is answer

KNN works very bad on

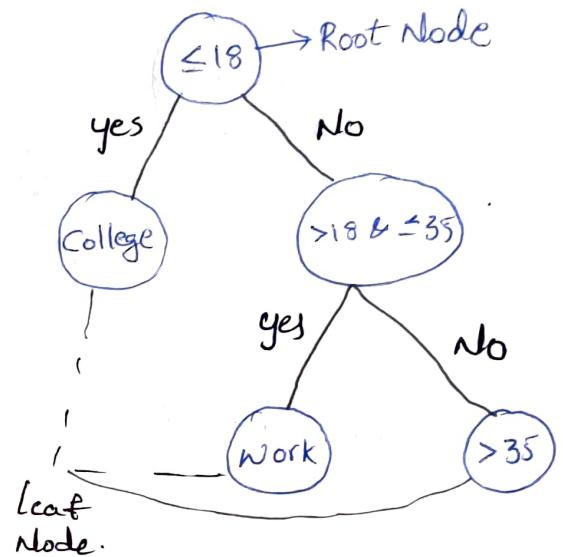
- (i) outlier
- (ii) Imbalanced dataset

# Day - 4 :-

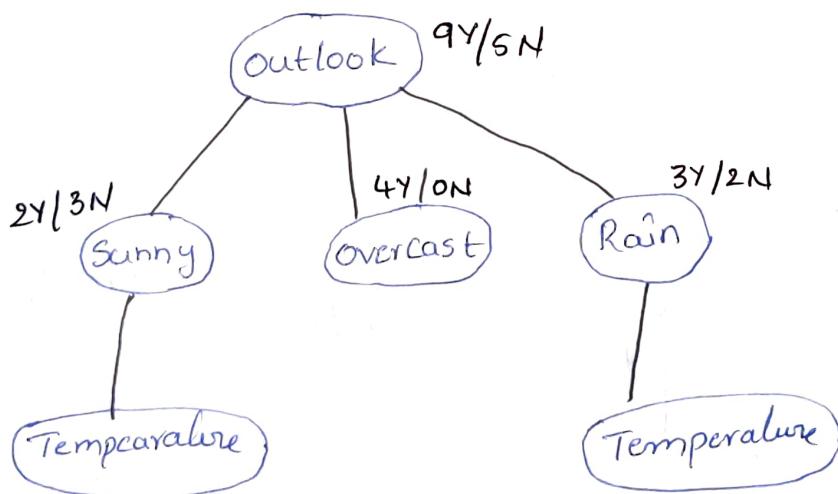
- ① Decision Tree classification
- ② Decision Tree Regression
- ③ Practical Implementation
- ④ Ensemble Techniques

## Decision Tree :-

```
if age ≤ 18 :  
    Print("college")  
  
elif (age > 18 & age ≤ 35) :  
    Print("work")  
  
else :  
    Print("Reactive")
```



## Dataset



Pure split  
Impure split

① Purity  $\Rightarrow$  Pure split ?

Entropy

Gini Impurity

Entropy  $E \{0, 1\}$

② How the features are selected ?  $\Rightarrow$  Information Gain

$p^+$  = probability of YES

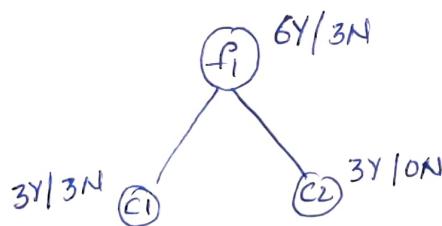
$p^-$  = probability of NO

Entropy

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

Gini Impurity

$$\begin{aligned} G.I &= 1 - \sum_{i=1}^n (P_i)^2 \\ &= 1 - [(P_+)^2 + (P_-)^2] \end{aligned}$$



$$\stackrel{C2}{=} -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$$= -1 \log_2 1 - 0$$

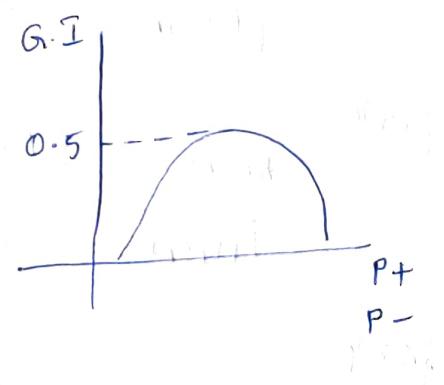
$$= 0 - 0$$

$= 0 \Rightarrow$  pure split

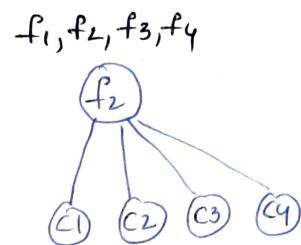
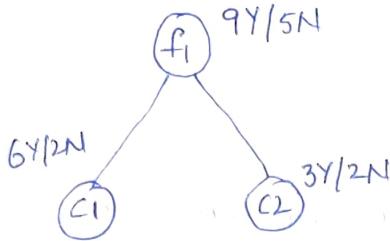
$$\begin{aligned} \stackrel{C1}{=} & -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\ & = 1 \Rightarrow \text{Impure split} \end{aligned}$$

$$\begin{aligned} f1 & 2Y/2N \\ G.I &= 1 - \left[ \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] \\ &= 1 - \left( \frac{1}{4} + \frac{1}{4} \right) \end{aligned}$$

$$\begin{aligned} &= 1 - \frac{1}{2} \\ &= 0.5 \end{aligned}$$



④ How to select which feature to split?



Information Gain

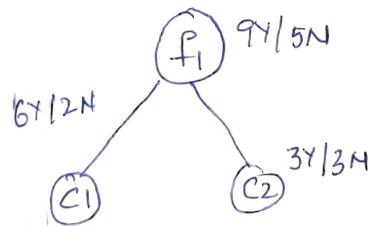
$$\text{Gain}(S, f_1) = H(S) - \sum_{r \in \text{val}} \frac{|S_r|}{S} H(S_r)$$

$H(S)$  = Entropy of feature

$|S_r|$  = Sample at category

$S$  = Total Sample

$H(S_r)$  = Entropy of category



{f1}

$$\begin{aligned} H(S) &= -P_+ \log_2 P_+ - P_- \log_2 P_- \\ &= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \end{aligned}$$

$$H(S_{f1}) = 0.94$$

{C1}

$$H(C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$H(C_1) = 0.81$$

{C2}

$$H(C_2) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$H(C_2) = 1$$

$$\text{Gain}(S, f_1) = H(S_{f_1}) - \sum_{\text{Eval}} \frac{|S_{rl}|}{S} H(S_{rl})$$

$$= 0.94 - \left[ \frac{8}{14} (0.81) + \frac{6}{14} (1) \right]$$

Gain( $S, f_1$ ) = 0.041

} which feature should start first  
if      Gain( $S, f_2$ ) = 0.051

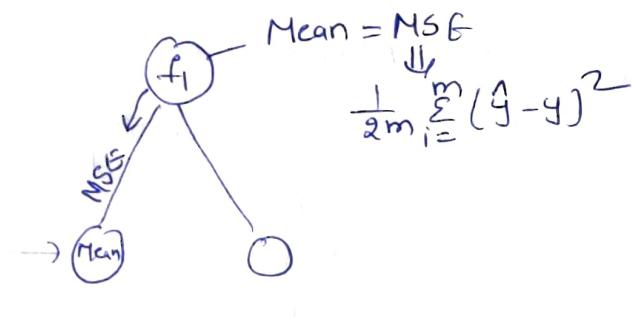
⊗ Node which has more gain starts first.

- Entropy takes more time to calculate algorithm
- Gini Impurity takes less time than entropy.

### Decision Tree Regressor $f$

$$f_1 \quad f_2 \quad \text{Op}$$

$$= \left\{ \begin{array}{l} \text{Mean} \\ \text{Min} \\ \text{Max} \end{array} \right.$$



### Hyperparameters $f$

⊗ Decision tree leads to overfit

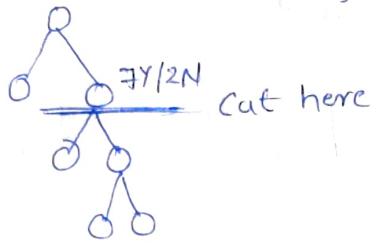
To avoid overfitting

(i) Post pruning

(ii) pre pruning

## Post pruning :-

Here, for yes  $\Rightarrow 70$ , so cut tree here

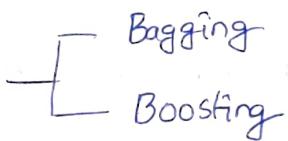


## Pre-pruning :-

max\_depth, max\_leaf  $\Rightarrow$  hyper parameters

## Day 5 :-

### ① Ensemble Techniques

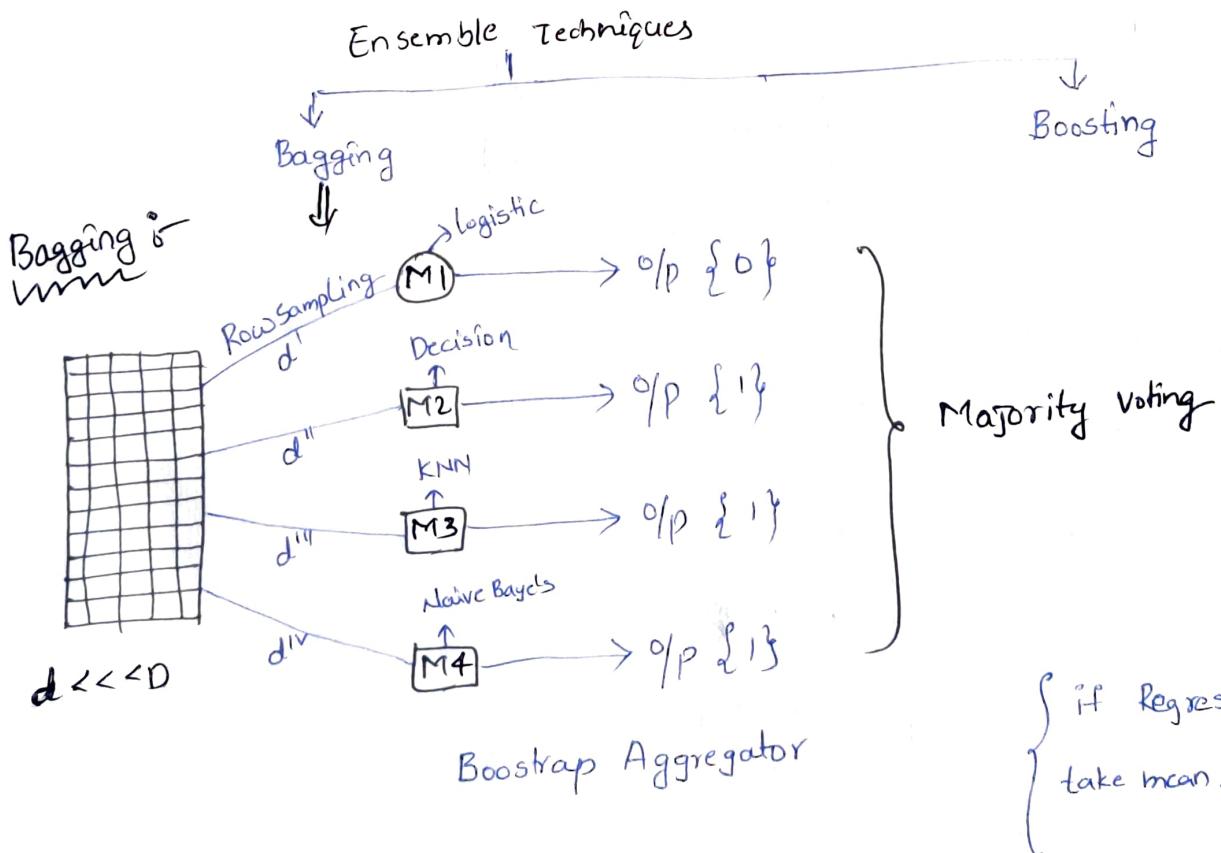


### ② Random Forest

### ③ AdaBoost

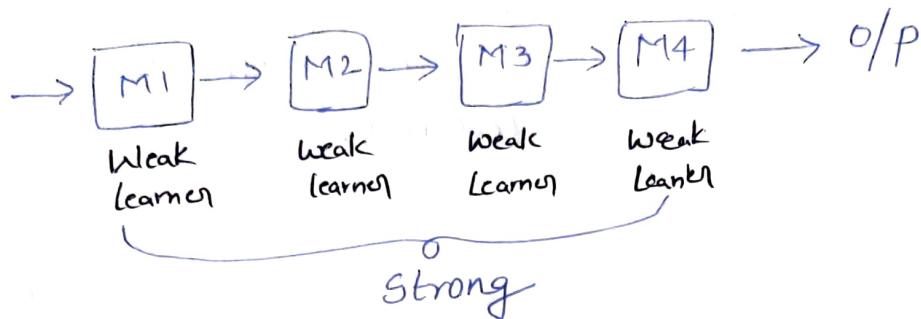
### ④ XGBoost

## Ensemble Techniques :-



## Boosting :-

Series of models



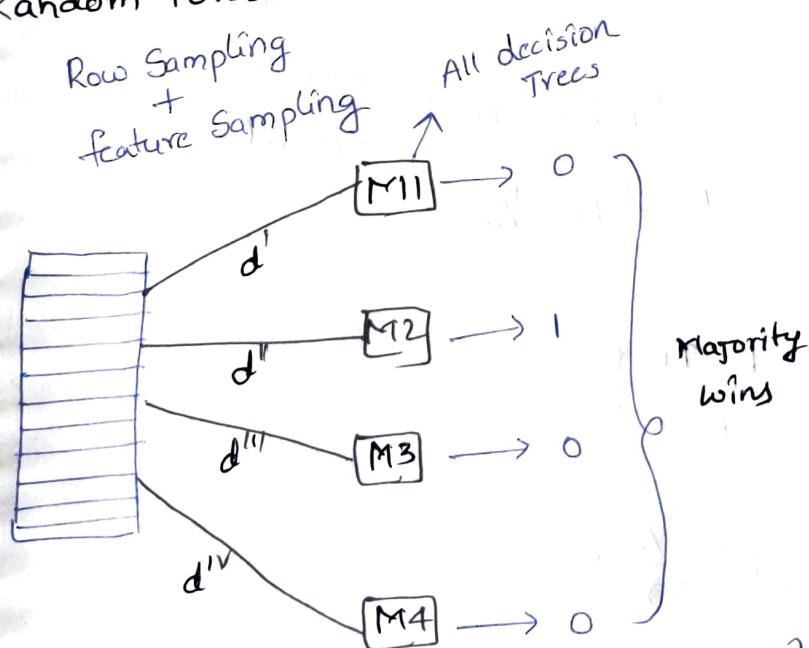
## BAGGING

- Random Forest classifier
- Random Forest Regressor

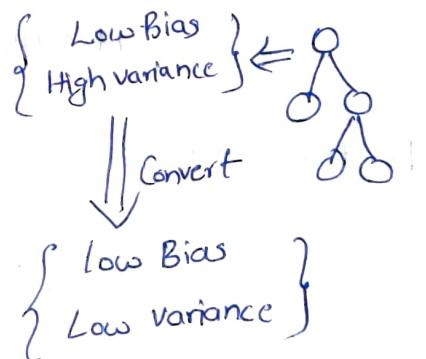
## BOOSTING

- AdBoost
- Gradient Boost
- XGBoost

## Random Forest classifier & Regression



What is problem of DT?  
↓  
overfitting



(i) Is normalization required?  
↓  
No? in decision Tree.

is normalization required in DT ?  $\Rightarrow$  No

is standardization required in KNN ?  $\Rightarrow$  Yes

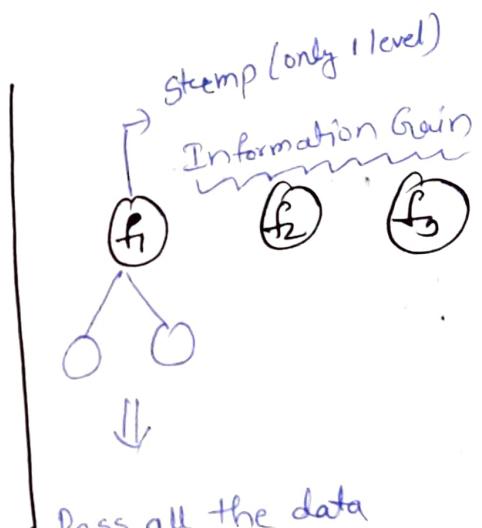
is random forest effect by outlier ?  $\Rightarrow$  No

$\hookrightarrow$  It gives equal opportunity to select features.

## Boosting

AdaBoost  $\leftarrow$  Sum of weight = 1

Total	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	$f_6$	$f_7$	Op	Weight
7	-	-	-	-	-	-	-	Y	$\frac{1}{7}$
	-	-	-	-	-	-	-	H	$\frac{1}{7}$
	-	-	-	-	-	-	-	N	$\frac{1}{7}$
	-	-	-	-	-	-	-	Y	$\frac{1}{7}$
	-	-	-	-	-	-	-	Y	$\frac{1}{7}$
	-	-	-	-	-	-	-	N	$\frac{1}{7}$



Pass all the data  
and find how many  
wrongly predicted

$$\text{Total Error} = TE = \frac{\text{wrongly}}{\text{Total}}$$

$$\text{if suppose wrong}=1 \Rightarrow TE=\frac{1}{7}$$

Performance of STUMP :

$$\boxed{\frac{1}{2} \log_e \frac{1-TE}{TE}}$$

$$P_S = \frac{1}{2} \log_e \left( \frac{1-\frac{1}{7}}{\frac{1}{7}} \right) \\ = 0.895$$

New Sample weight

$$\text{Correct} = \text{weight} * e^{-P_S}$$

$$\frac{1}{7} * e^{-0.895} = 0.05$$

$$\text{Incorrect} = \text{weight} * e^{P_S}$$

$$\Rightarrow \frac{1}{7} * e^{0.895} = 0.349$$

New weight

$$\begin{array}{r}
 0.05 \\
 0.05 \\
 0.05 \\
 0.349 \\
 0.05 \\
 0.05 \\
 0.05 \\
 0.05 \\
 \hline
 0.649
 \end{array}$$

Normalized weight

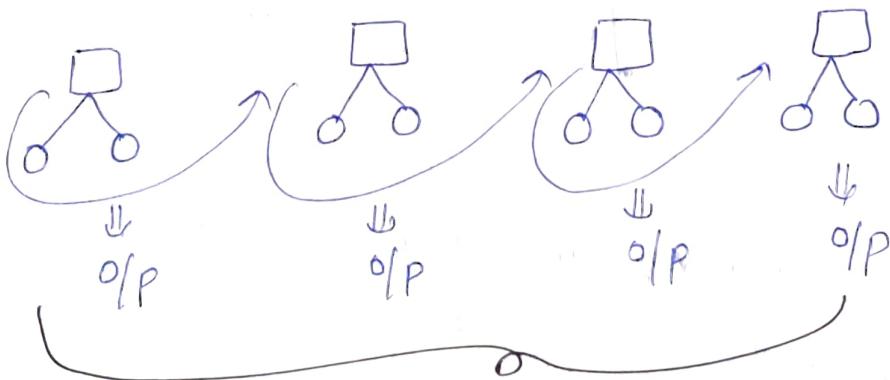
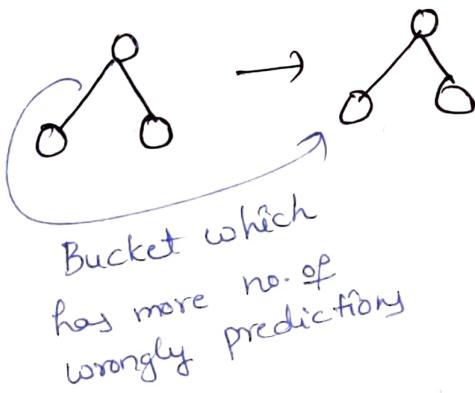
$$\begin{array}{r}
 0.07 \\
 0.07 \\
 0.07 \\
 0.537 \\
 0.07 \\
 0.07 \\
 0.07 \\
 \end{array}$$

Buckets

[0 - 0.07]
[0.07 - 0.14]
[0.14 - 0.21]
[0.21 - 0.749]
[0.749 - 0.751]
[ - ]
[ - ]

Here, sum of weight  $\neq 1$   
So, Normalize by divide 0.649

This having high  
wrong values  
↓  
Passed to next node



Majority  $\rightarrow$  classification  
Mean  $\rightarrow$  Regression

Linear Regression  $\rightarrow$  White Box Model

Random Forest  $\rightarrow$  Black Box model

Decision Tree  $\rightarrow$  White Box model

ANN  $\rightarrow$  Black Box model

## Day-6 b

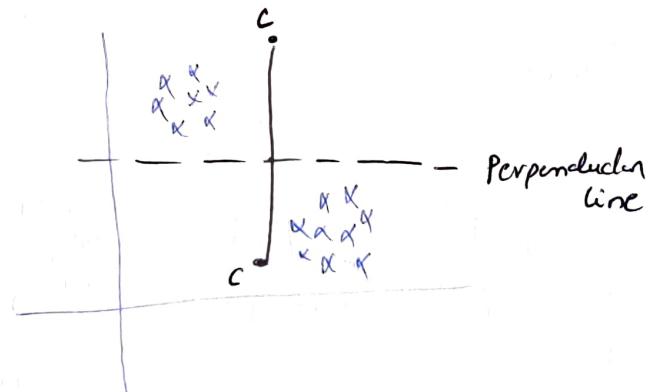
### Unsupervised ML

- ① KMeans clustering
- ② Hierarchical clustering
- ③ Silhouette score
- ④ DBSCAN clustering

K-Means clustering :-

k-centroids

- (i) we try with different k-values
- (ii) Initialize k number of centroids
- (iii) Compute avg to update centroids

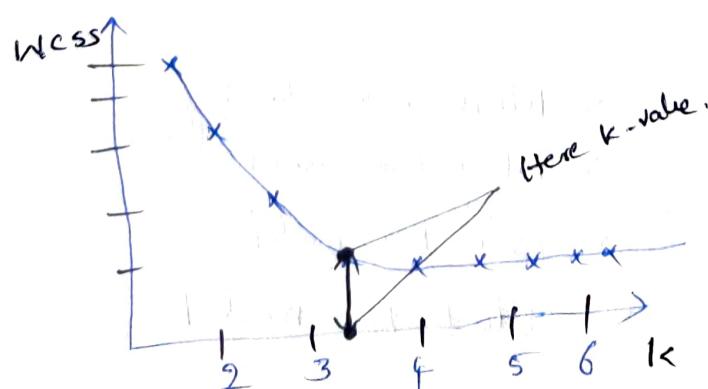


Note :- Initialize centroid very far  $\Rightarrow$  K-Means ++

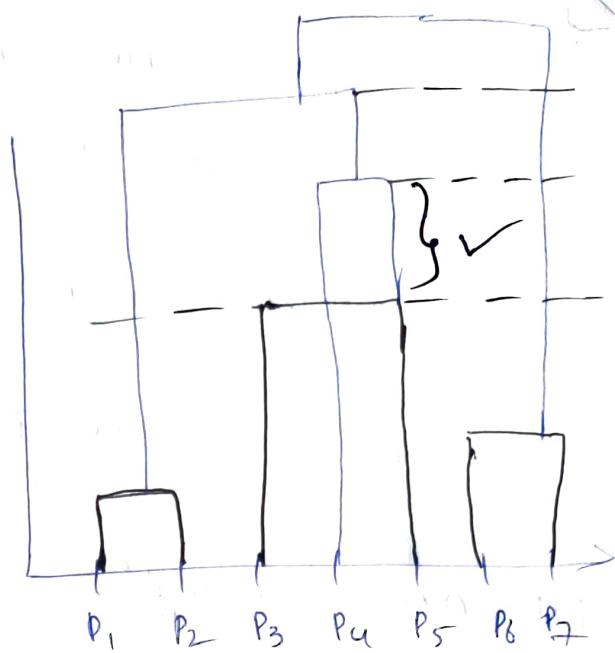
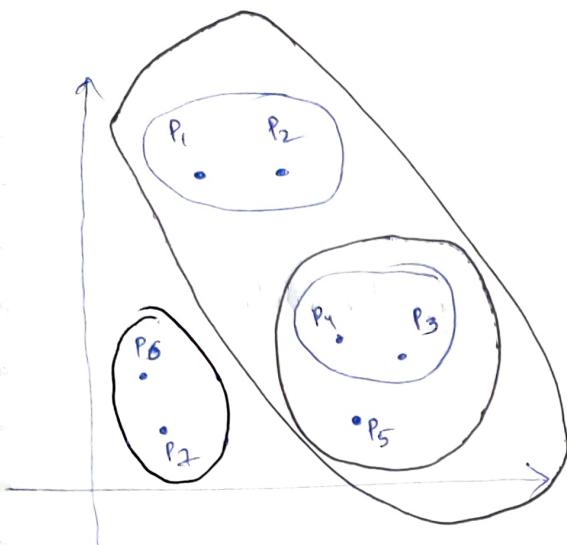
How to decide k-value ?

Elbow method

WCSS - Within cluster sum of squares



## ② Hierarchical clustering



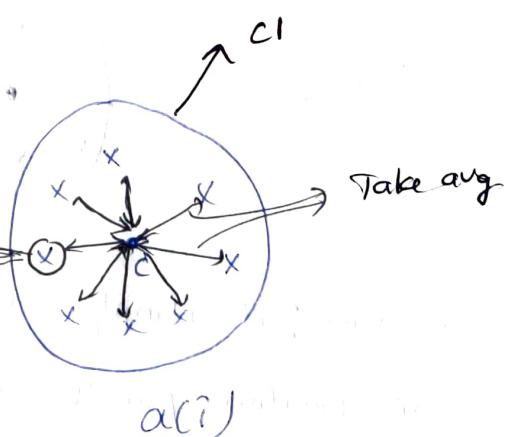
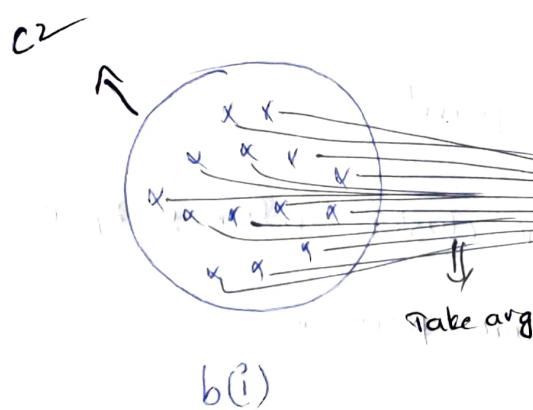
Dendrogram

- \* Find the largest vertical line that has no horizontal line pass through.

→ Maximum time taken by k-Means or Hierarchical clustering?

Sol: Hierarchical clustering

Silhouette Score(clustering) :-



$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \left\{ \begin{array}{l} \text{if } |c(i)| > 1 \end{array} \right.$$

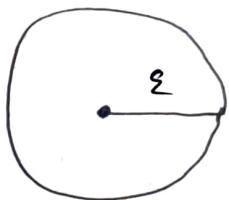
values = {-1 to 1}

if  $S(i)$  is towards +1 then it is more good model.

## DBSCAN clustering b'

Density Based spatial clustering of Applications with Noise

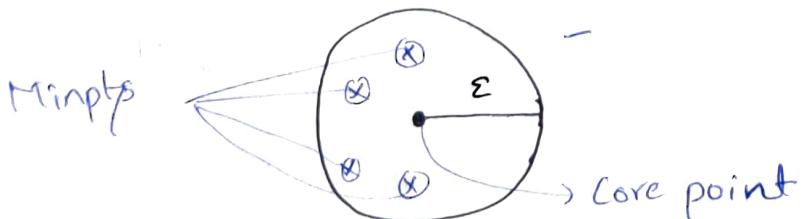
- i) Min points
- ii) Core points
- iii) Border points
- iv) Noise points
- v) Epsilon ( $\epsilon$ )



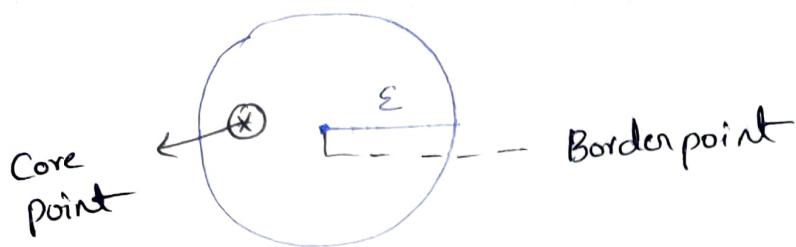
\* DBSCAN used to avoid outlier or noise points.

\* Minpts is hyperparameter.

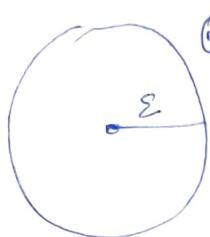
if  $\text{minpts} = 4$



Condition: if  $\text{minpts}$  are inside the cluster radius of Epsilon the epsilon points called as core point.



\*) within only one core point, then it becomes border point.



① No points inside

② This is outside point, neglect  
→ It is Noise point

### Bias :-

It is phenomenon that skews the result of an algorithm in favour or against an idea.

### Variance &

The variance refers to the changes to the model when using different portions of training or test data.

### Day - 7 G

① Xg boost classifier

② Xg boost Regressor

③ SVM

④ SVR

# Xgboost classifier:- Extreme Gradient Boosting

<u>Salary</u>	<u>Credit</u>	<u>Approval</u>	<u>{Residual}</u> $(1 - pr)$
$\leq 50$	B	0	-0.5
$\leq 50$	G	1	0.5
$\leq 50$	G	1	0.5
$\leq 50$	G	0	-0.5
$\geq 50$	G	1	0.5
$> 50$	N	1	0.5
$> 50k$	N	0	0.5

↓                          ↑

Base Model  $\Rightarrow \text{pr}(= 0.5)$

① Create a binary decision tree using feature

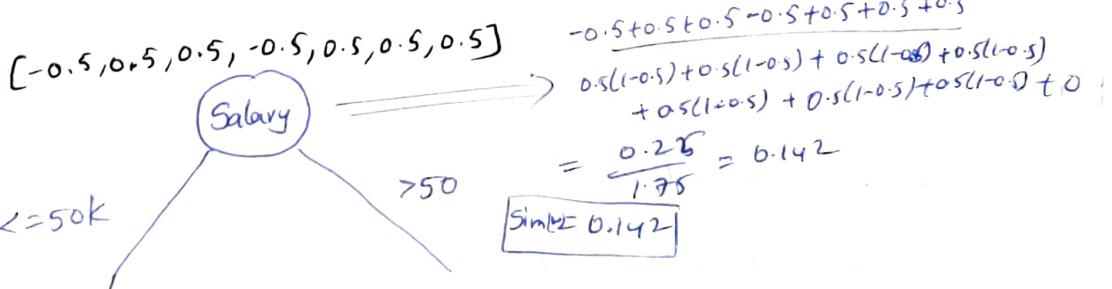
② Calculate Similarity weight

$\lambda$  = hyperparameter

$$= \frac{\sum (\text{Residual})^2}{\sum (\text{pr}(1-\text{pr})) + \lambda}$$

③ Information Gain

Solt



Assume  
 $\lambda = 0$

$$\text{Sim}(w) = \frac{(-0.5 + 0.5 + 0.5 - 0.5)^2}{(0.5)(1-0.5) + (0.5)(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} = \frac{0}{4} = 0$$

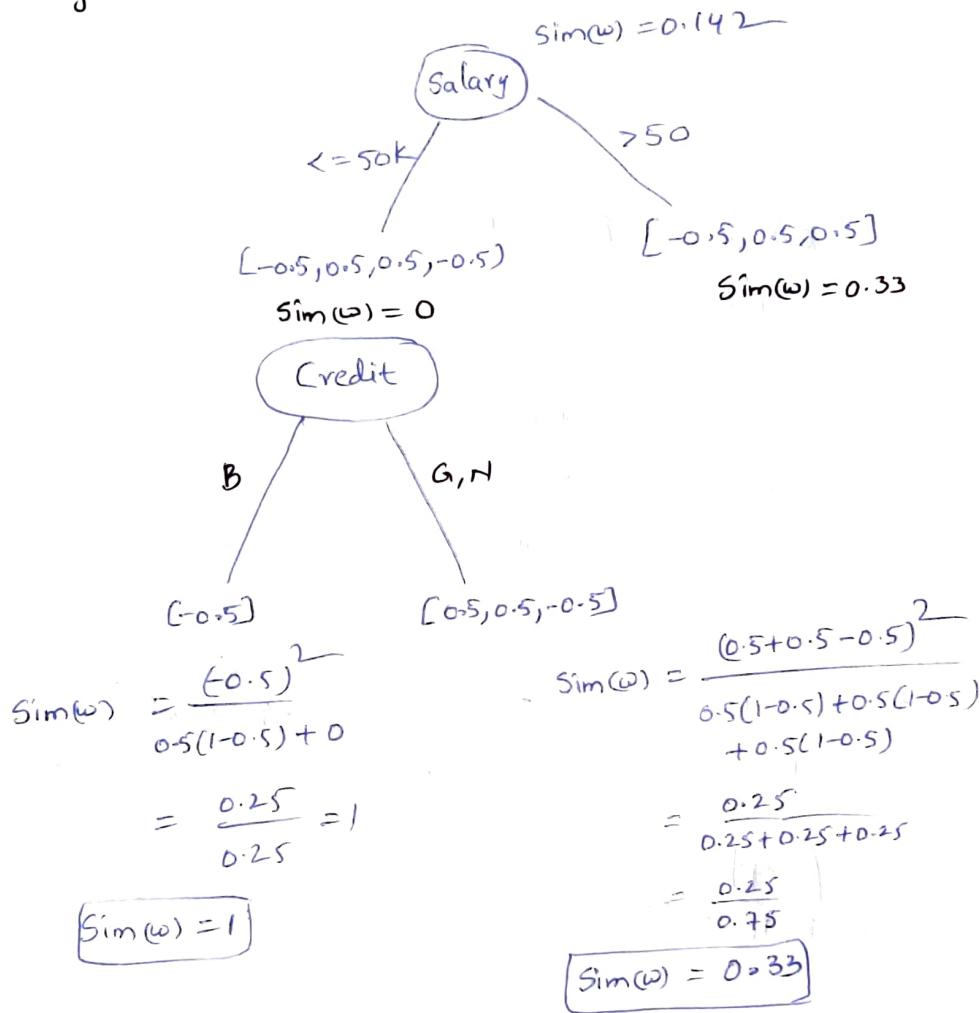
$\text{Sim}(w) = 0$

$$\text{Sim}(w) = \frac{(-0.5 + 0.5 + 0.5)^2}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} = \frac{0.25}{1.75} = 0.33$$

$\text{Sim}(w) = 0.33$

$$\begin{aligned}
 \text{Information Gain} &= (\text{Sum of } \text{sim}(\omega) \text{ of branch nodes}) - \text{sim}(\omega) \text{ of Root Node} \\
 &= 0 + 0.33 - 0.14 \\
 &= 0.19
 \end{aligned}$$

Divide again above DT



$$G \cdot I = 1 + 0.33 - 0$$

$$(G \cdot I = 2.33)$$

$$\text{for base model } 1 \quad \log \left( \frac{P}{P-P} \right) = \log \frac{0.5}{1-0.5} = 0$$

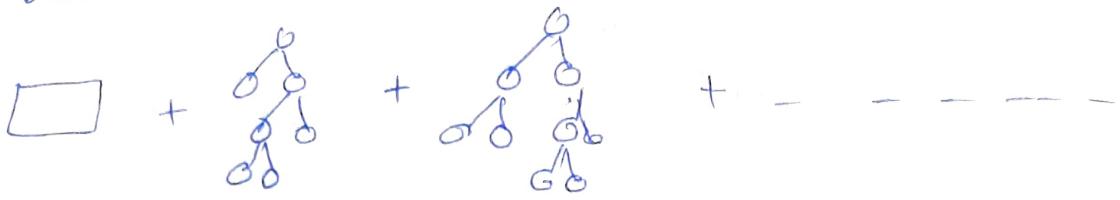
To find real probability

$$1 - [0 + \alpha (\text{sim}(\omega) \text{ of leaf Node})]$$

$$1 - [0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n)]$$

\* Xgboost is also blackbox model.

$$\text{Base} + \alpha_1(DT_1) + \alpha_2(DT_2) + \dots + \alpha_n(DT_n)$$



XG Regressor :-

Exp	Gap	Salary	Residual
2	Yes	40K	-11K
2.5	Yes	42K	-9K
3	No	52K	+1K
4	No	60K	+9K
4.5	Yes	62K	+11K

$$\boxed{\text{Base}} \rightarrow \text{avg of o/p} = 51K$$

$$\text{Sim}(w) = \frac{\sum (\text{Residuals})^2}{\text{No. of residuals} + 1}$$

$$\text{Sim}(w) = \frac{(-11+9+1+9+11)^2}{5+0}$$

$$\text{Sim}(w) = \frac{121}{1+0}$$

$$\boxed{\text{Sim}(w) = 121}$$

$$\begin{array}{c} \boxed{\text{Exp}} \\ \leftarrow [ -11, 9, 1, 9, 11 ] \\ \leq 2 \quad \quad \quad > 2 \end{array}$$

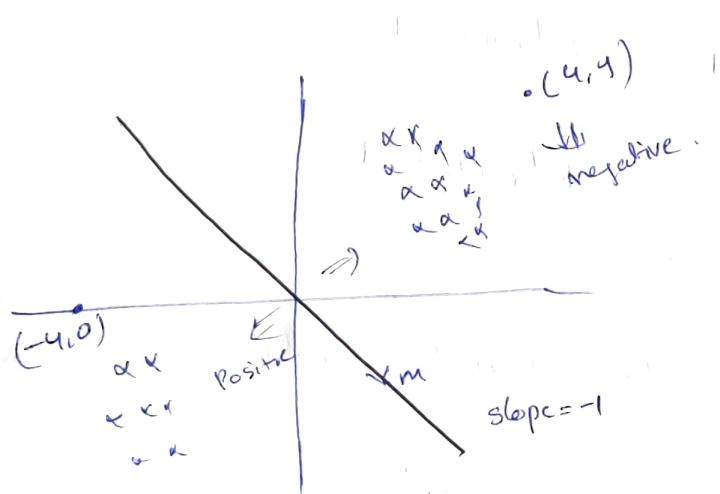
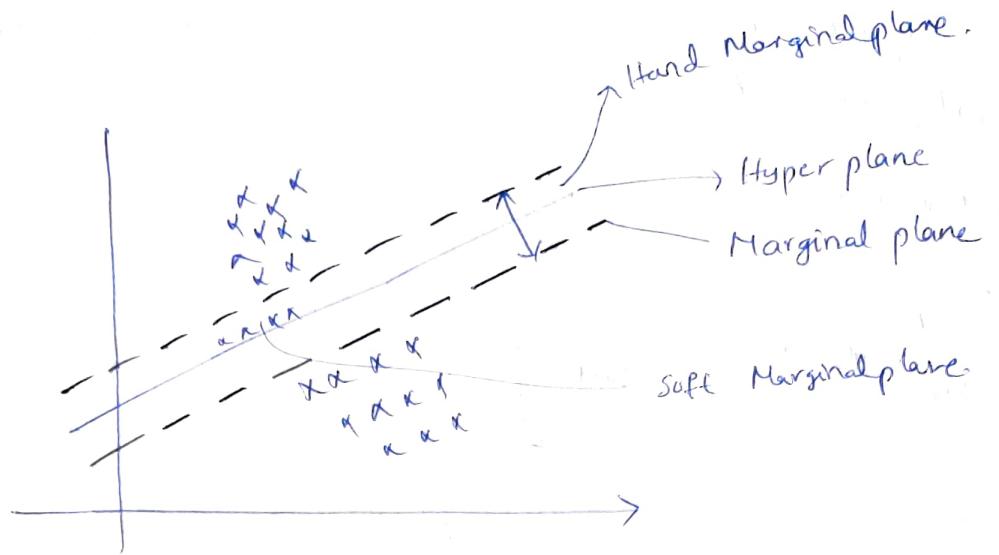
$$\begin{aligned} \text{Sim}(w) &= \frac{(-9+1+9+11)^2}{4+0} \\ &= \frac{144}{4} \end{aligned}$$

$$\boxed{\text{Sim}(w) = 36}$$

$$\text{Information Gain} = \frac{121 + 36 - 0.2}{156.8}$$

$$= 156.8$$

## SVM :-



$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} [-4, 0]$$

$$= 4$$

$\Rightarrow +ve$  value

$$y = mx + c$$

$$y = w_1 x_1 + w_2 x_2 + \dots + b$$

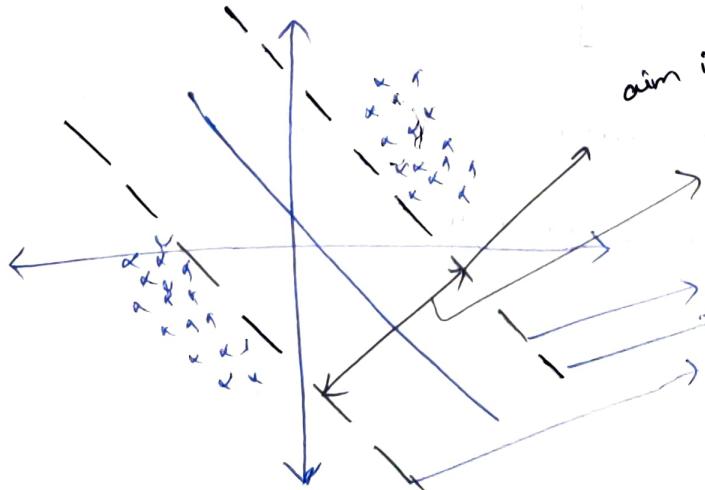
$$y = \mathbf{w}^T \mathbf{x} + b$$

$$(4, 4)$$

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} [4, 4]$$

$$= -4$$

-ve value



aim is to increase or maximize distance

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 2$$

$$\mathbf{w}^T \mathbf{x} + b = 1$$

$$\mathbf{w}^T \mathbf{x} + b = 0$$

$$\mathbf{w}^T \mathbf{x} + b = -1$$

$$\begin{aligned} \omega^T x_1 + b &= 1 \\ \omega^T x_2 + b &= -1 \\ \hline \omega^T(x_1 - x_2) &= 2 \end{aligned}$$

divide by  $|\omega|$  to remove magnitude.

$$\boxed{\frac{\omega^T}{|\omega|} (x_1 - x_2) = \frac{2}{|\omega|}}$$

Aim is to maximize  $\frac{2}{|\omega|}$

such that

$$y_i \begin{cases} +1 & \omega^T x_i + b \geq 1 \\ -1 & \omega^T x_i + b \leq -1 \end{cases}$$

Major aim

$$y_i * (\omega^T x_i + b) \geq 1 \quad \text{for correct points}$$

$$\boxed{\text{Maximize}_{(\omega, b)} \frac{2}{|\omega|} \Leftrightarrow \text{Min}_{(\omega, b)} \frac{|\omega|}{2}}$$

$$\boxed{\text{Min}_{(\omega, b)} \frac{|\omega|}{2} + \sum_{i=1}^m \xi_i}$$

$\xi_i$  = summation of the distance of the wrong data points

$c_i$  = How many errors we have