

# ML ALGORITHM

22/12/24

1  
10/12/24

Day-1

\* Linear Regression

⇒ 1. Introduction to ML (AI vs ML vs DL vs DS)

⇒ 2. Supervised ML and Unsupervised ML

⇒ 3. Linear Regression (Maths & Geometric Intuition)

⇒ 4.  $R^2$  & Adjusted  $R^2$

⇒ 5. Ridge and Lasso Regression

⇒

1. AI vs ML vs DL vs DataScience

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



Netflix → Action → Recommendation

↳ Comedy → Recommendation

Amazon → iPhone → Recommend to (MacBookPro)

→ ML → Subset of AI

→ Start tools to analyze, visualize, predict, forecast

→ ML is a subset of ML

→ mimic a human brain

→ multilayer neural network

→ DataScience is a part of Everything.

vision

27/12/24

DATA SCIENCE SEM

2

## Machine Learning & Deep Learning

Supervised ML vs Un-Supervised ML

Regression Classification

Clustering Dimensional Reduction

→ Supervised machine learning is by ML

Input → We have done dataset information DA ←

Age	weight
24	62
25	63
27	72
27	62
29	73

Age → Hypothesis → O/P  
Weighted

→ Independent Feature → Age  
→ Dependent Feature → weight

vision

3  
→ Supervised machine Learning is

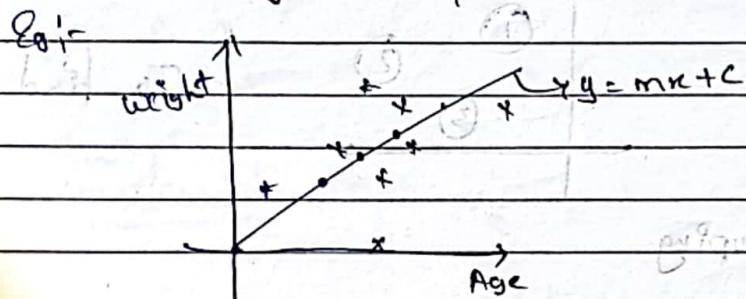
## 1) Regression Problem Statement

Age	weight	→ OIP Variable
24	7.2	
23	7.1	
25	7.5	

Continuous Variable

Regression Problem Statement

→ Whenever your output is continuous it is a  
Regression Problem Statement.



## 2) Classification Problem Statement

Ques -

OIP

No. of hours No. of play hrs No. of sleep hrs

PIF

Play -

P

Not Play -

F

Sleep -

P

Not Sleep -

F

→ Whenever you have in your output fix number of categories that become a classification problem

Categ yes / No

Pass / Fail

Spam / not spam

→ 2 outputs become a "Binary Classification"

→ more than 2 outputs "Multiclass Classification"

## X Un-Supervised Machine Learning

### i) Clustering

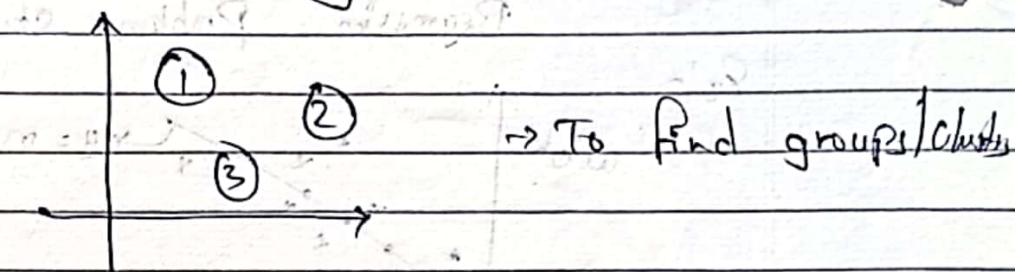
#### ii) Dimensional Reduction

↳ Clustering is

Salary Age

$\left\{ \begin{array}{l} \text{No dep variable} \\ \text{No dependent variable} \end{array} \right\}$

"Clustering"  $\rightarrow$  Customer Segmentation



$\Rightarrow$  Grouping

Eg:-

From launching our "product"

$\hookrightarrow$  Product 1  $\rightarrow$  Rich people

Product 2  $\rightarrow$  Middle class

Product 3  $\rightarrow$  Poor people

"Add Marketing"

"Customer Segmentation"

### iii) Dimensional Reduction

$\rightarrow$  Lower Dimension

$$\begin{matrix} 1000 \\ 100 \\ 100 \end{matrix} \rightarrow \begin{matrix} 100 \end{matrix}$$

↓ ↓ ↓

Similarly it is possible through using "PCA", etc.

"LDA"

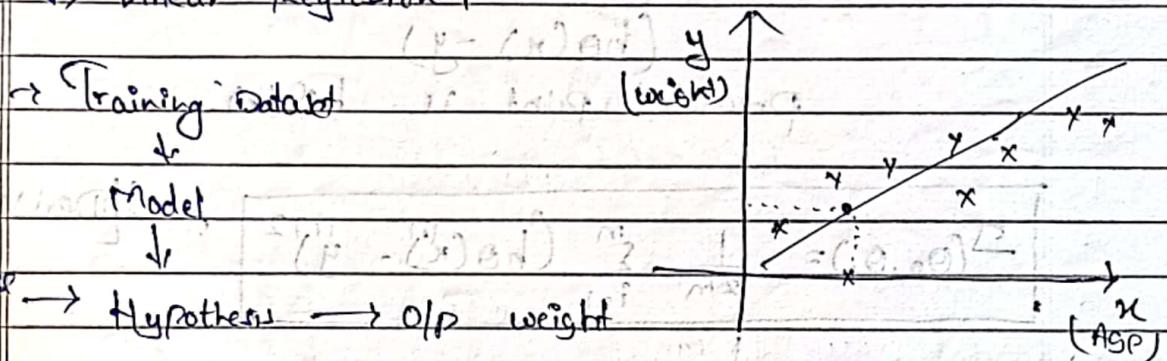
## Supervised

Un-supervised

- 1) Linear Regression
- 2) Ridge & Lasso
- 3) Logistic Regression
- 4) Decision Tree
- 5) Ada Boost
- 6) Random Forest
- 7) Gradient Boosting
- 8) K-boost
- 9) k-means
- 10) DB Scan
- 11) Hierarchical clustering
- 12) k-Means + Nearest neighbor clustering
- 13) PCA
- 14) LDA

## Explanation :-

### 1) Linear Regression :-



→  $y$  is a linear function of  $x$ .

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

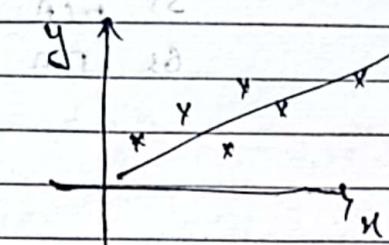
$$\left. \begin{array}{l} y = \theta_0 + \theta_1 x \\ \text{how: } \theta_0 + \theta_1 x \end{array} \right\}$$

Equation of straight line :-

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

$$y = \theta_0 + \theta_1 x^{(i)}$$

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



$\theta_0$  = Intercept (when  $x=0$ )

$\theta_1$  = slope ( $b$ ) coefficient

$x_i$  = data points

→ Cost Function

⇒ distance formula ⇒ b/w predicted &

$(h_{\theta}(x_i) - y_i)$   
predicted point is  $h_{\theta}(x_i)$ .

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

[Derivation]

Squared Error Function

⇒ In this we need to solve :-

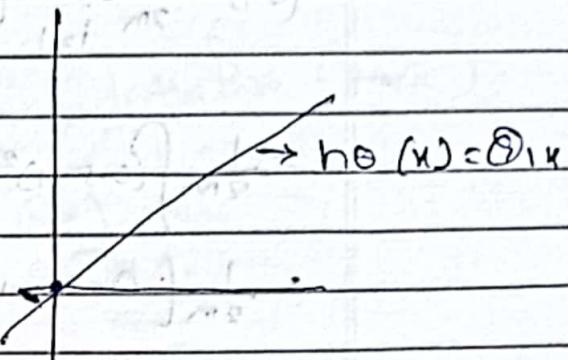
$$\begin{aligned} \text{minimize}_{\theta} \quad & \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ \text{Subject to} \quad & \downarrow \\ \text{minimize}_{\theta_0, \theta_1} \quad & J(\theta_0, \theta_1) \end{aligned}$$

\* Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \text{If } \theta_0 = 0$$

\* Example :-

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

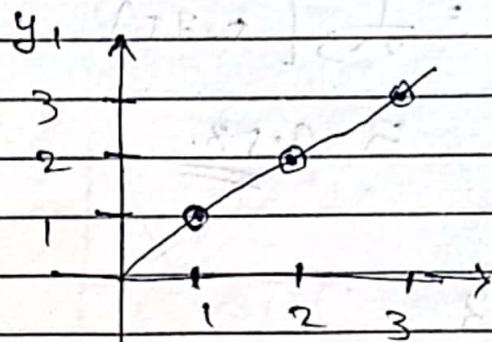


$$\theta_1 = 1$$

$$(1, 1)$$

$$(2, 2)$$

$$(3, 3)$$



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

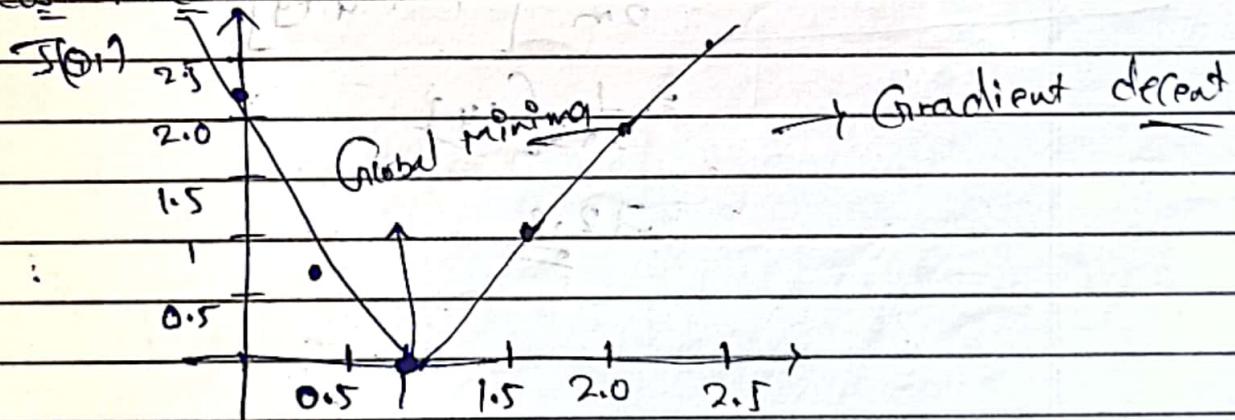
$$= \frac{1}{2m} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$$J(\theta_1) = 0$$

$$\theta_1 = 1$$

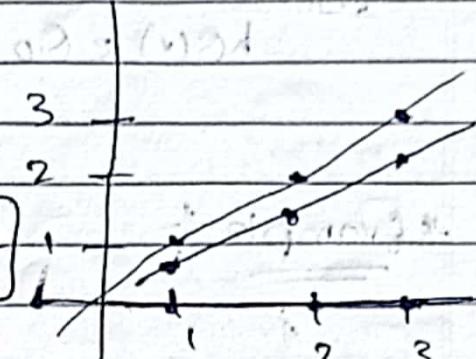
$$J(\theta_1) = 0$$

\* Cost Function



Now,  $Q_1 = 0.5$

$$J(Q_1) = \frac{1}{2m} \sum_{i=1}^3 (h_\theta(x^{(i)}) - y^{(i)})^2$$



$$= \frac{1}{2m} \left[ (0.5 - 1)^2 + (1.5 - 2)^2 + (2.5 - 3)^2 \right]$$

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

$$= \frac{1}{3} [5.75]$$

$$\approx 0.58$$

3

Now  $Q_1 = 0$

when  $Q_1 = 0$  it will be come in minimum

4

Now  $Q_1 = 0$

$$J(Q_1) = \frac{1}{2m} \left[ (0-1)^2 + (0-2)^2 + (0-3)^2 \right]$$

$$= \frac{1}{2m} [1 + 4 + 9]$$

$$= \frac{1}{3} [14]$$

$$\approx 2.3$$

1.0 0.5 2.5 2.0

a

## Conversions      Algorithms

Repeat until Convergence.

$$\left. \frac{\partial}{\partial \theta_j} I(\theta_0, \theta_1) \right|_{\theta = \theta^*}$$

Example

$\Rightarrow$  Gradient Descent Algorithmen

Repeat until Convergence

$$Q_j := Q_j - \lambda \frac{\partial}{\partial q_j} J(Q_0, \psi_1)$$

I is always 0 and 1

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

$$\stackrel{i=0}{\rightarrow} \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

## Conversion algorithm

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

Convergence Algorithm

10

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

Repeat until Converge

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

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

\* Performance metrics

$R^2$  and Adjusted  $R^2$

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

residual :- Sum of Squared residual

Total :- Total sum of Squares ranges 0 and 1  
Small number → big number  
High → small number

⇒ Adjusted  $R^2$  formula

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

N = No of data points

P = No of Predictions

⇒ If P is increasing  $(N-P-1)$  value is decreasing

big No sum value to get  
small No

Let's do this

$R^2$  is useful for understanding how much variance is explained by your model

Adjusted  $R^2$  is better for model comparison and ensuring you don't overfit by adding irrelevant predictors

adjusted  $R^2$

adjusted  $R^2$

(less well)

adjusted  $R^2$  - Now writing down

(less well)

adjusted  $R^2$

(less well)

adjusted  $R^2$

adjusted  $R^2$  - less well if removing terms

D

adjusted  $R^2$  - less well if removing terms

D

adjusted  $R^2$  - less well if removing terms

D

adjusted  $R^2$  - less well if removing terms

D

adjusted  $R^2$  - less well if removing terms

D

adjusted  $R^2$  - less well if removing terms

D

vision

ability

less well

less well

## \* Agenda

1. Ridge & Lasso Regression
2. Assumption of Linear Regression
3. Logistic Regression
4. Confusion matrix
5. Practical for linear, Ridge, Lasso, & logistic

Explanation :-

### Overfitting

(Low Bias)

Model performs well  $\rightarrow$  Training data

Fails to perform well  $\rightarrow$  Test Data

(High Variance)

### Underfitting

(High Bias, High variance)

Model Accuracy is bad with training data

Model Accuracy is also bad with test data.

! -

### Model - 1

Training ACC = 90%.

Test ACC = 80%.



Overfitting

Low Bias  
High Variance

### Model - 2

Training ACC = 92%.

Test ACC = 91%.



Generalized model  
Low Bias  
High Variance

### Model - 3

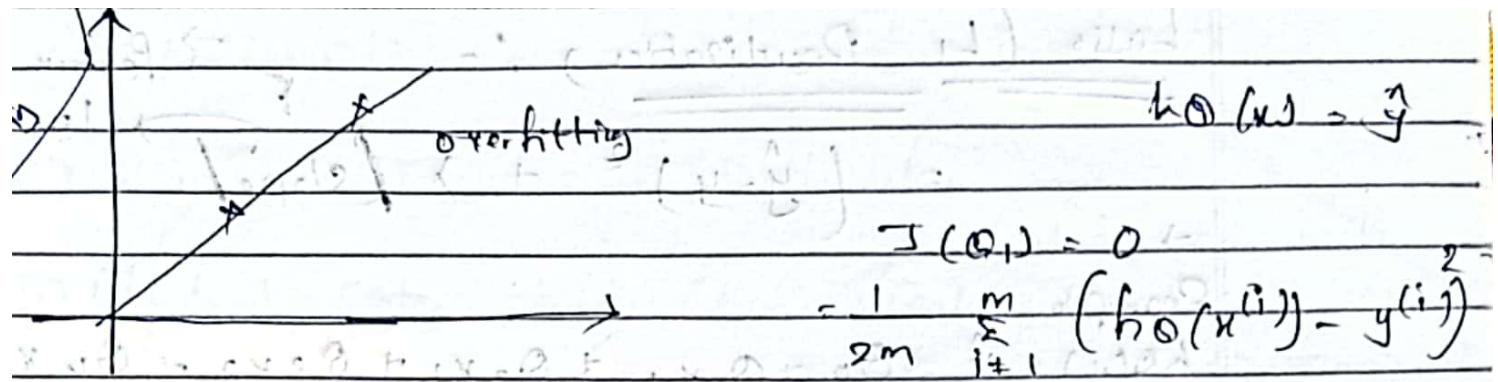
Training ACC = 70%.

Test ACC = 65%.



Underfitting

High Bias  
High Variance



here,  $J(\theta_0, \theta_1) = (y^{(i)} - \hat{y}^{(i)})^2 + \lambda(\text{slope})^2$

$= 0$

Ridge ( $L_2$  Regularization)

when  $\lambda$  value is 4

" $\lambda$ " is hyper parameter

Abp. when  $\theta_0 = 0$ ,

$$h_0(x) = \theta_1 x$$

$\downarrow$  Slope

$$\theta_1 + 1(2)^2 = 5 \downarrow \text{minimum}$$

$$0 + 1(2)^2$$

$$= (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope}_c)^2$$

$$= (\text{small value}) + 1(1.25)^2$$

$$= (\text{small value}) + 2.25$$

$$= 1 \approx 3$$

$\Rightarrow$  iteration {hyper parameter}

$R^2$ , Adjusted  $R^2$

$\Rightarrow$  that is lambda ( $\lambda$ )

$\Rightarrow$  Hyper parameter

$\Rightarrow$  we get cross validation

It will use

Lasso ( $L_1$  Regularization) :-

→ Feature selection

$$| \alpha_1 + \alpha_2 + \alpha_3 + \dots + \alpha_n |$$

$$= (\hat{y} - y)^2 + \lambda / \text{slope}$$

Example 6

$$h(x) = \hat{y} = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \alpha_4 x_4 + \dots + \alpha_n x_n$$

For

① preventing overfitting }  $L_1$  Regularization

② feature selection }  $L_2$  Regularization

(minimizing)

→ Cross validation

•  $\hat{y}(x) = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n$

→

\* Ridge Regression ( $L_2$  Normalization)

$$\text{Cost Function} = (h(x^{(i)}) - y^{(i)})^2 + \lambda / \text{slope}^2$$

1) Purpose: Preventing Overfitting

\* Lasso Regression ( $L_1$  Reg)

$$\text{Cost Function} = (h(x^{(i)}) - y^{(i)})^2 + \lambda / \text{slope}$$

\* Purpose: 1) prevent overfitting

2) Feature Selection

$$(h(x)) \rightarrow \text{minimize } |\alpha_1 + \alpha_2 + \alpha_3 + \dots + \alpha_n|$$

## \* Assumption of linear Regression :-

- ① R in Normal / Gaussian Distribution  $\rightarrow$  model will get trained well
- ② Standard Scaler / standardization: (Scaling data)  $\rightarrow$  Z-Score  $\text{mean} - c, \text{stand deviation}(a) = z$
- ③ Linearity
- ④ Multi-collinearity

$$\Rightarrow x_1, x_2, \dots, y$$

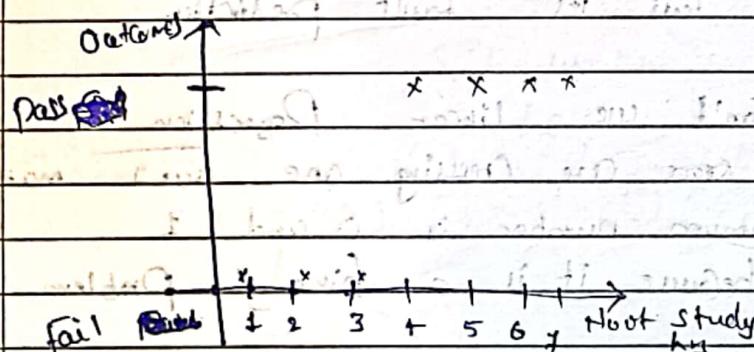
Transformation In Plateau Factor?

## \* Logistic Regression (Classification) :-

$\rightarrow$  Binary Classification

No. of study      No. of play      Pass / Fail

-	-	P
-	-	F
-	-	P or F



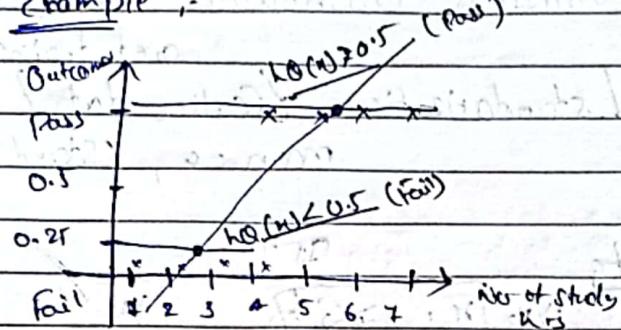
$\therefore$  If study less than 3 hrs "Fail"

$\therefore$  If study more than 3 hrs "pass"

$\Rightarrow$  Can we solve "linear Regression"?

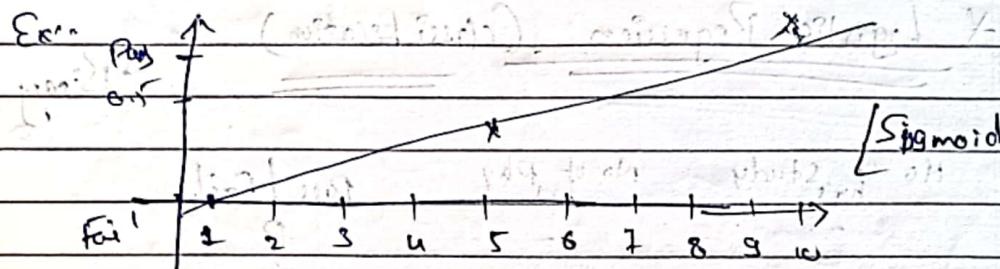
Yes we can solve...

Example :-



$$\begin{aligned} h\theta(x) > 0.5 \Rightarrow 0 \rightarrow \text{Fail} \\ h\theta(x) \leq 0.5 \Rightarrow 1 \rightarrow \text{Pass} \end{aligned}$$

\* but we don't use because same confusion let's check



$$h\theta(x) = 0.10$$

$\Rightarrow$  here 5 is pass but it's look like a fail its fault prediction

$\Rightarrow$  Can't use linear Regression

we are crossing one our maxima & minimum number is 0 and 1  
because it is a binary problem

## \* Decision Boundary logistic Regression :-

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

$$h_{\theta}(x) = \sigma \underset{\text{Transforms}}{\underset{x}{\text{linear}}} \rightarrow h_{\theta}(x)$$

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

(Squash linear straight line)

hypothesis

log

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1)$$

$$\text{Let } z = \theta_0 + \theta_1 x_1$$

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

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

This function is called  
Sigmoid (or) Logistic Function

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}}$$

It will work well

because

Ex: graph :- It is being to squash the function

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

0.5

If without confusion

not get false value

$$\text{When } g(z) \geq 0.5 \quad z \geq 0.5$$

$\Rightarrow$  Training set

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

$$y \rightarrow 0 \text{ or } 1 \xrightarrow{\text{feed }} \text{OLD}$$

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

$$z = \theta_0 + \theta_1 x_1$$

$$z = \theta_1 x$$

Change parameter  $\theta_1$

$\Rightarrow$  Cost Function

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

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

here,

Logistic Regression Cost Function:

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

$\Rightarrow$  We can't use this cost function

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

because  $\frac{1}{1 + e^{-(\theta_1 x)}}$  this non-convex function

$\Rightarrow$  Initialize it, non-convex function

\* Non-Convex Function

$\rightarrow$  This is related to Gradient descent

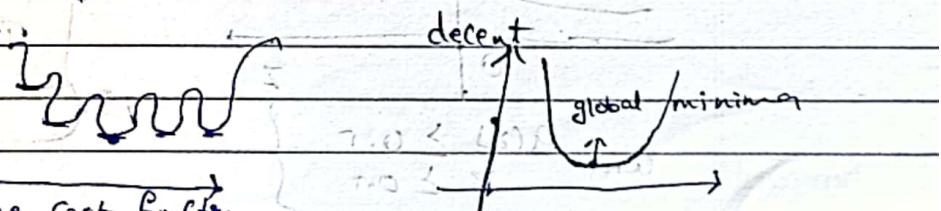
\* Convex Function

$\rightarrow$  This is related to Gradient descent

$\rightarrow$  There are so many local minima

So we can't reach global

$\rightarrow$  vision minima never that is the reason we can't ... real life

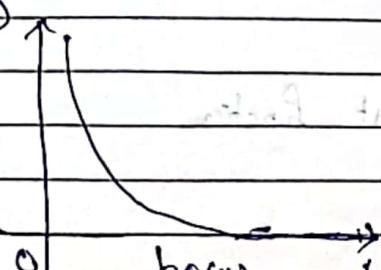


## Logistic Regression Cost Function

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

$\rightarrow$  if  $y=1$ ,  $J(\theta)$  is small  
 $\rightarrow$  if  $y=0$ ,  $J(\theta)$  is large

$J(\theta)$

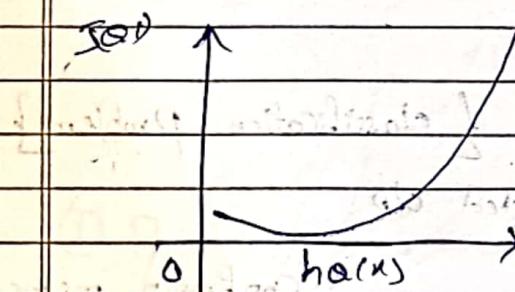


Cost = 0 if  $[y=1, h_{\theta}(x)=1]$

Classification problems are often

$\rightarrow$  if  $y=0$

$J(\theta)$



$\rightarrow$

$$\text{Cost}(h_{\theta}(x^{(i)}), y) = \begin{cases} -\log(h_{\theta}(x^{(i)})) & \text{if } y=1 \\ -\log(1-h_{\theta}(x^{(i)})) & \text{if } y=0 \end{cases}$$

Final cost function

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

vision

$$\left. \begin{array}{l} \text{if } y=1 \\ C_1(\text{h}_\theta(x^i); y) = -\log(h_\theta(x^i)) \\ \text{if } y=0 \\ C_0(\text{h}_\theta(x^i); y) = -\log(1-h_\theta(x^i)) \end{array} \right\} \begin{array}{l} \text{both functions} \\ \text{are proved} \\ \text{by} \\ \text{for Cost function} \end{array}$$

$$\Rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m [C_1(\text{h}_\theta(x^i); y^i) + C_0(\text{h}_\theta(x^i); y^i)]$$

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

→ This becomes entire cost function

Repeat until Convergence

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

\*

→ Performance metrics {classification problem}

$x_1$	$x_2$	$y$	$\hat{y}$	actual O/P	Predicted O/P	Confusion matrix	Calculation
-	-	0	1	-	-	-	-
-	-	1	1	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-	-	-
-	-	0	1	-	-	-	-
-	-	1	0	-	-	-	-
-	-	0	0	-	-</		

		Actual		
		1	0	
Predicted	1	TP	FP	Confusion matrix
	0	FN	TN	This always gives right OLP

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

$$= \frac{3+1}{3+2+1+1} = \frac{4}{7}$$

$$= 0.57 \Rightarrow 57\%$$

→ When we have 0's are 900 } It is imbalance data  
 1's are 100 }

When,

$$\begin{cases} 0's \rightarrow 600 \\ 1's \rightarrow 400 \end{cases} \quad \text{Balanced data}$$

⇒

$$\begin{cases} 0 : 900 \\ 1 : 100 \end{cases} \quad \text{but It is a Bias data}$$

$$\left\{ \text{Model} \rightarrow 0 = \frac{900}{1000} = 80\% \right\}$$

① Precision

$$\frac{TP}{TP+FP}$$

formula's

② Recall

$$\frac{TP}{TP+FN}$$

③ F-Score

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

$$\beta+1 = (1+\beta) \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$= \frac{2(\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$

$$\frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

vision

Harmonic mean

$$\frac{2xy}{x+y}$$

This is called a "Harmonic mean" thus focuses both  
False negative & False Positive

$$\beta = 0.5 \quad F_{0.5} \text{ score} = \frac{(1 + (0.5)^2) P \times R}{(0.25) P + R}$$

$$\beta = 2 \quad FN > FP$$

F. 2 Score

$$(P \times R) / (P + R)$$

$$(True Positive \times True Positive) / (True Positive + False Positive)$$

$$(TP \times TP) / (TP + FP)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

$$(TP \times TP) / (TP + FP + FN)$$

06/01/25  
NS

## Day-3

23

### \* Today Agenda :-

- 1) Practical Simpler Examples
- 2) Naive Baye's Intuition
- 3) kNN algorithms

- 1. Linear Regression
- 2. Ridge & Lasso
- 3. Logistic Regression

⇒ Linear Regression, Ridge and Lasso :-

### Here Code :-

```
→ Importing all libraries
from sklearn.datasets import fetch_california_housing
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

→ # Load the data-set
df = fetch_california_housing()

→ We check type here of dataset
type(df)

→ We print dataset
df      # Now we shown some data & target data

# We convert dataset to DataFrame
dataset = pd.DataFrame(df.data)          → In dataset one column name
dataset.columns = df.feature_names
dataset.head()   # Print a tiny data set in DataFrame
                ↳ This data is independent Feature.
print(dataset.head()) # It print full
```

## Dividing the dataset into independent and dependent features.

$X = \text{dataset.iloc[:, :-1]}$  # skipping the last feature  
 # Independent Feature

$y = \text{dataset.iloc[:, -1]}$  # dependent feature

$X.head()$

# Now we use algorithms

\* Linear Regression

from sklearn.linear\_model import LinearRegression

from sklearn.model\_selection import cross\_val\_score

lin\_reg = LinearRegression() # (cross\_val\_score)

mse = cross\_val\_score(lin\_reg, X, y) # (cross\_val\_score)

, Scoring = 'neg\_mean\_squared'

Error', cv=5)

(1) original - size 10 | 1 2 3 4 5 6 7 8 9 10

train data | test data | train data

(2) random - size 10 | 1 2 3 4 5 6 7 8 9 10

train data | test data | train data

test data

mean\_mse =

mean\_mse = np.mean(mse)

print(mean\_mse) # output :- -0.5582901717186536

C:\Users\91750\OneDrive\Desktop\ML\Linear Regression

1. Importing the libraries

2. Importing the dataset

3. Preprocessing the dataset

4. Splitting the dataset into training set and test set

### \* Ridge algorithms

```

from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV
ridge = Ridge()
params = {'alpha': [1e-15, 1e-10, 1e-8, 1e-3,
                   1e-2, 1, 5, 10, 20]}
ridge_regressor = GridSearchCV(ridge, params, scoring
                               = 'neg_mean_squared_error', cv=5)
ridge_regressor.fit(x, y)
print(ridge_regressor.best_params_)
print(ridge_regressor.best_score_)
# output:
{'alpha': 20}
-0.5581020035625638

```

# Ridge tries to reduce the overfitting

### \* Lasso Regression

```

from sklearn.linear_model import Lasso
from sklearn.model_selection import GridSearchCV
lasso = Lasso()
params = {'alpha': [1e-15, 1e-10, 1e-8, 1e-3,
                   1e-2, 1, 5, 10, 20]}
lasso_regressor = GridSearchCV(lasso, params, scoring
                               = 'neg_mean_squared_error', cv=5)
lasso_regressor.fit(x, y)
print(lasso_regressor.best_params_)
print(lasso_regressor.best_score_)

```

\* Using train test split

```
X = dataset.iloc[:, :-1] # Independent feature
y = dataset.iloc[:, -1] # dependent feature
```

```
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import cross_val_score
lin_reg = LinearRegression()
lin_reg.fit(X_train, y_train)
mse = cross_val_score(lin_reg, X_train, y_train, scoring
                      = 'neg_mean_squared_error', cv=5)
print(mean(mse))
print(mse)
```

⇒ Ridge Regression Algorithms :-

```
from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV
ridge = Ridge()
params = {'alpha': [1e-15, 1e-10, 1e-8, 1e-3, 1e-2, 2,
                   1, 10, 15, 20]}
```

```
ridge_regressor = GridSearchCV(ridge, params, scoring='neg_mean_squared_error', cv=5)
```

```
ridge_regressor.fit(X_train, y_train)
```

```
print(ridge_regressor.best_params_)
```

```
print(ridge_regressor.best_score_)
```

## ~~\* Lasso Regression~~

```
from sklearn.linear_model import Lasso  
from sklearn.model_selection import GridSearchCV  
lasso = Lasso()  
params = {'alpha': [1e-15, 1e-10, 1e-8, 1e-3, 1e-2,  
                  1, 2, 10, 20]}
```

Lasso regressor = GridSearchCV (lasso, param\_grid, scoring = "neg\_mean\_squared\_error", cv=5)

lasso-regressor · fit(X\_train, y\_train)

Print Class - regressor . best params -)

```
print(classifier.best_score_)
```

→ # Using R<sup>2</sup> Score (R-Square)

`y_pred = lasso_regressor.predict(X_test)`

```
from sklearn.metrics import r2_score
```

$$H2\_Score1 = \text{np.sum}(Score(Cy\_pred, Cy\_test1))$$

Print (re-scored) : taken = return(4-3) in loop (a)

## \* Logistic Regression :-

```
from sklearn.linear_model import LogisticRegression
import pandas as pd
from sklearn.datasets import load_breast_cancer
df = load_breast_cancer()
df = df.data, df['target'] # independent variable
X = pd.DataFrame(df[0], columns=df[1])
X.head()
y = pd.DataFrame(df[1], columns=['Target'])
y['Target'].value_counts() # using train-test-split:
```

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)
params = [{"C": [1, 5, 10], "max_iter": [100, 150]}]
model1 = LogisticRegression(C=10, max_iter=100)
GridSearchCV(Estimator=model1, param=grid, params,
            scoring='f1', cv=5)
```

```
model1.fit(X_train, y_train)
```

```
model1.best_params_
```

```
model1.best_score_
```

```
y_pred = model1.predict(X_test)
```

```
y_pred
```

```
from sklearn.metrics import Confusion_matrix, classification_report
```

```
Confusion_matrix(y_test, y_pred)
```

```
Accuracy_Score(y_test, y_pred)
```

```
Print(classification_report(y_test, y_pred))
```

## \* Naive Bayes

Intuition :-

{specifically use for classification}



## {Base Theory}

Rolling a Dice

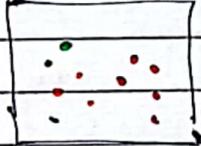
$$\{1, 2, 3, 4, 5, 6\}$$

$$\left. \begin{array}{l} P(1) = \frac{1}{6}, P(2) = \frac{1}{6} \\ P(3) = \frac{1}{6} \end{array} \right\}$$

{Independent Events}

Dependent Event :-

I have a Bag of marble



$$P(R) = \frac{5}{11}, P(G) = \frac{3}{10}$$

← First Element removed

← Second Element removed

$$P(R \text{ and } G) = P(R) * P(G|R)$$

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

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

{we can write like this  
"Yes"}

Formula :-

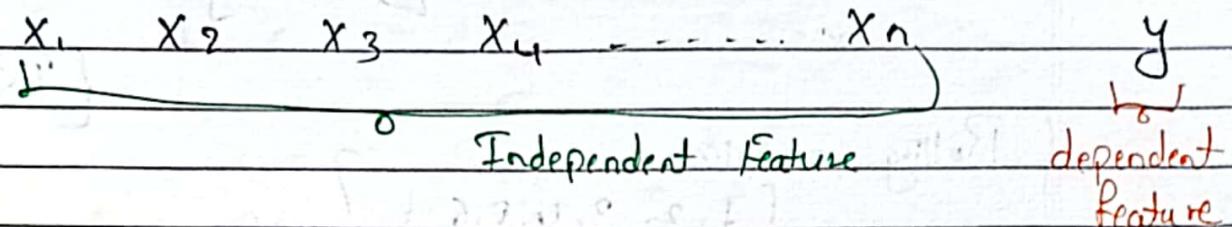
$$P(A) * P(A|B) = P(B) * P(A|B)$$

$$P(A) = \frac{P(B) * P(A|B)}{P(A)}$$

→ Bayes

Theorem  
Crux Naive Bayes

→ I have Some Features like



$$P(y | x_1, x_2, x_3, \dots, x_n) = P(y) * P(x_1, x_2, \dots, x_n | y)$$

$$P(x_1, x_2, x_3, \dots, x_n)$$

$$= P(y) * P(x_1 | y) * P(x_2 | y) * P(x_3 | y) * \dots * P(x_n | y)$$

$$P(x_1) * P(x_2) * P(x_3) * \dots * P(x_n)$$

⇒ Data set is

$$x_1 \quad x_2 \quad x_3 \quad x_4 \quad y$$

$$\begin{array}{cccc} \text{Remember} & \text{Condition 1: } & \text{Condition 2: } & \text{Condition 3: } \\ \text{Data 1: } & (A, b, 1, \text{yes}) & (A, b, 0, \text{no}) & (A, b, 1, \text{no}) \\ \text{Data 2: } & (A, b, 0, \text{no}) & (A, b, 1, \text{no}) & (A, b, 0, \text{yes}) \end{array}$$

$$P(y = \text{yes} | x_i) = P(\text{yes}) * P(x_1 | \text{yes}) * P(x_2 | \text{yes}) * P(x_3 | \text{yes}) * P(x_4 | \text{yes})$$

$$\text{remember: } \text{Constant (Fixed)} \rightarrow P(x_1) * P(x_2) * P(x_3) * P(x_4)$$

$$P(y = \text{no} | x_i) = P(\text{no}) * P(x_1 | \text{no}) * P(x_2 | \text{no}) * P(x_3 | \text{no}) * P(x_4 | \text{no})$$

$$\text{Constant (Fixed)} \rightarrow P(x_1) * P(x_2) * P(x_3) * P(x_4)$$

$x_i$        $\overset{\text{yes}}{\underset{\text{no}}{\text{I}}}$

$$P(\text{yes} | x_i) = 0.13$$

$$\geq 0.5 \Rightarrow 1$$

$$< 0.5 \Rightarrow 0$$

$$P(\text{no} | x_i) = 0.05$$

$$P(\text{yes}/\text{ni}) = \frac{0.13}{0.13+0.05} = 0.72 = 72\%$$

Final Answer

$$P(\text{No}/\text{ni}) = 1 - 0.72 = 0.28 = 28\%.$$

Independent Future

\* Dayset :-

	Day	Outlook	Temperature	Humidity	Wind	Play	Tenn
D <sub>1</sub>	Sunny	Hot	High	Weak	W	Y	
D <sub>2</sub>	Sunny	Hot	High	Strong	W	Y	
D <sub>3</sub>	Overcast	Hot	High	Weak	S	N	
D <sub>4</sub>	Rain	Mild	High	W	S	N	
D <sub>5</sub>	Rain	Cool	Normal	W	S	N	
D <sub>6</sub>	Rain	Cool	Normal	S	S	N	
D <sub>7</sub>	Overcast	Cool	Normal	S	S	N	
D <sub>8</sub>	Sunny	Mild	High	W	Y		
D <sub>9</sub>	Sunny	Cool	Normal	W	Y		
D <sub>10</sub>	Rain	Mild	Normal	W	Y		
D <sub>11</sub>	Sunny	Mild	Normal	S	Y		
D <sub>12</sub>	Overcast	Mild	Normal	S	Y		
D <sub>13</sub>	Overcast	Hot	High	W	N		
D <sub>14</sub>	Rain	Mild	High	S	N		

=> Outlook :-

	Probability(Yes)		Probability(No)	
	P(Y)	P(N)	P(Y)	P(N)
Sunny	2	3	2/5	3/5
Overcast	4	0	4/5	0/5
Rain	3	2	3/5	2/5
Total	9	5		

80 80

\* Temperature :-

	Yes	No	$P(Y)$	$P(N)$
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cold	3	1	3/9	1/5

Total      9      5

\* Play Tennis :-

	$P(\text{Yes})$	$P(\text{No})$
Yes	9	
No	5	
Total	14	

\* Play Tennis

	$P(\text{Yes})$	$P(\text{No})$
Yes	9	9/14
No	5	5/14
Total	14	

Now,

→ Test (Sunny, Hot) → o/P ?

$$\rightarrow P(\text{yes}, (\text{Sunny}, \text{Hot})) = P(\text{yes}) \times P(\text{Sunny} | \text{yes}) \times P(\text{Hot} | \text{yes})$$

$$= \frac{9}{14} \times \frac{2}{9} \times \frac{2}{5}$$

$$= \frac{2}{7} \times \frac{2}{9} = \frac{2}{63} = \underline{\underline{0.031}}$$

$$P(\text{No}) \times P(\text{Sunny}|\text{No}) \times P(\text{Hot}|\text{No})$$

$$\cancel{P(\text{No})} \times \cancel{P(\text{Sunny}|\text{No})} \times \cancel{P(\text{Hot}|\text{No})}$$

$$P(\text{No} | \text{Sunny, Hot}) = \frac{3}{14} \times \frac{2}{7} \times \frac{3}{5} \xrightarrow{\text{constant}} \text{constant}$$

$$\frac{3}{35} = \underline{\underline{0.085}}$$

$$\Rightarrow P(\text{yes} | \text{Sunny, Hot}) = 1 - 0.085 = 0.915 = 91.5\%$$

$$\Rightarrow P(\text{No} | \text{Sunny, Hot}) = 0.085 \quad \begin{cases} \text{Normal} \\ \text{2nd row} \end{cases}$$

$$0.085 / (0.085 + 0.915) = \underline{\underline{0.085}}$$

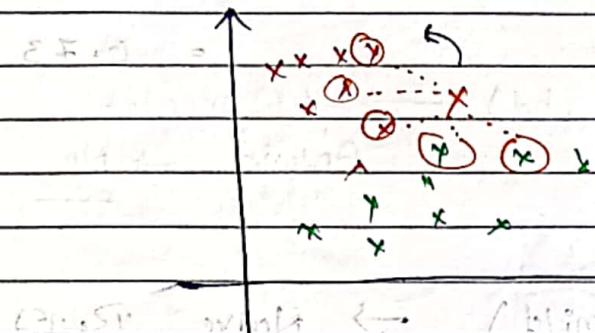
Letten  
 $\rightarrow (\text{Sunny, hot}) \rightarrow \text{Yes or No}$   
 Answer  $\rightarrow \underline{\text{No}}$

Assignment  
 Now, new sum:  
 $(\text{overcast, mild}) \rightarrow \text{Naive Bayes?}$

## 2) KNN Algorithm :- (K-Nearest Neighbour)

Classification → Regression

Classification → K-Nearest Neighbour



With  $K=3$

Calculate distance Red - 3

Green - 2

→ we use tree distance

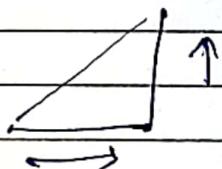
i) Euclidean Distance

ii) Manhattan Distance

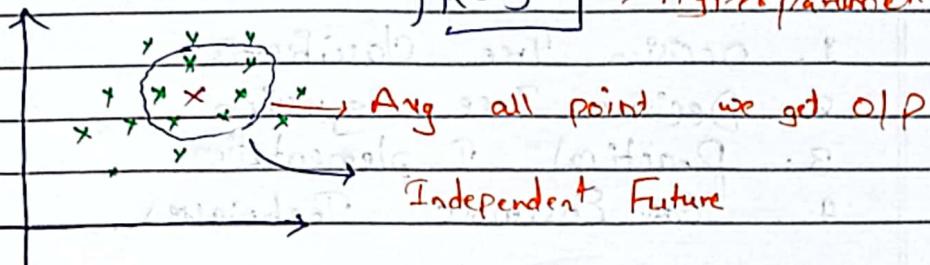
⇒ i) Euclidean Distance

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

ii) Manhattan Distance



## ii) Regression :-

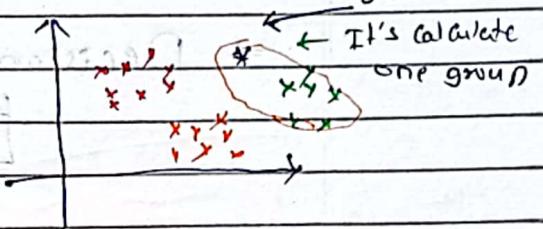
 $K=5$  $\rightarrow$  Hyperparameter

- \* Imp note
- \*  $\star$

- $\star$  KNN works very badly

- $\star$  Outliers

- $\star$  Imbalance of



REVIEW

Classification and Regression

(Classification)

QUESTION

Classification

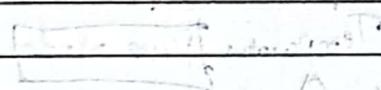
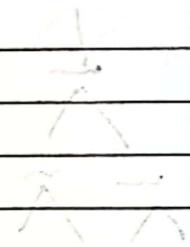
(Classification)

Classification  $\rightarrow$  (if we want) prediction

solution

Help

Help



A house is divided into two parts by a line.

## \* Agenda

1. Decision Tree Classification
2. Decision Tree Regression
3. Practical Implementation
4. Ensemble Techniques.

⇒ Decision Tree { Solving many usecases }

{ Regression }  
Classification

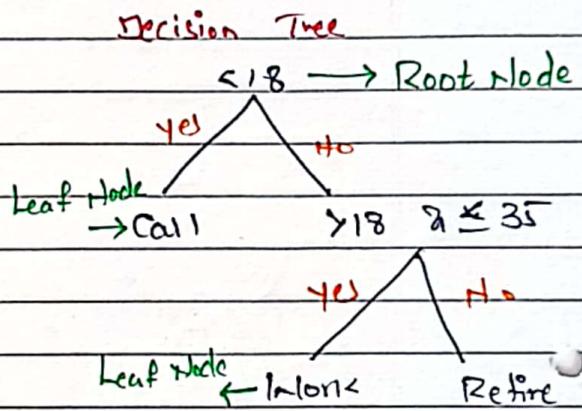
⇒

if age < 18:  
print ("Go to college")

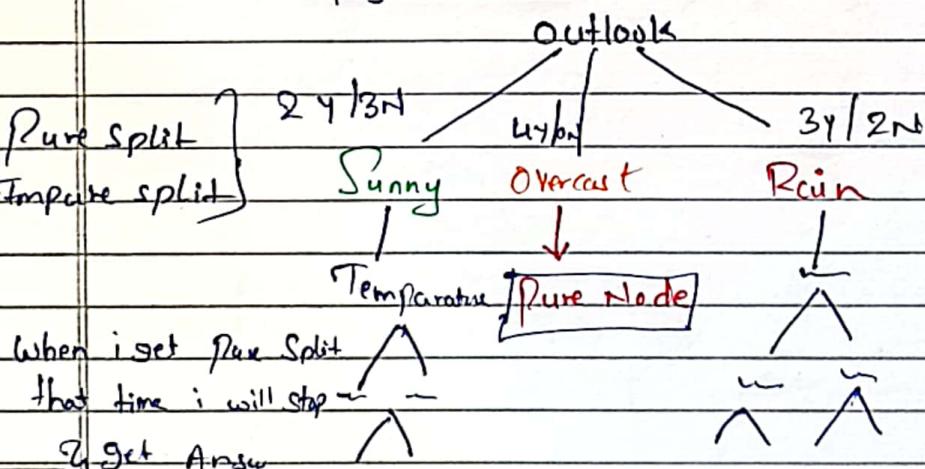
Elif age ≥ 18 and age ≤ 35:  
print ("Work")

Else:

print ("Retire")



## \* Decision (page No - 31) → Classification



1. Purity  $\rightarrow$  pure split ???

$\hookrightarrow$  Entropy

$\hookrightarrow$  Gini coefficient / Impurity

2. How the Features are selected

$\hookrightarrow$  Information gain ??

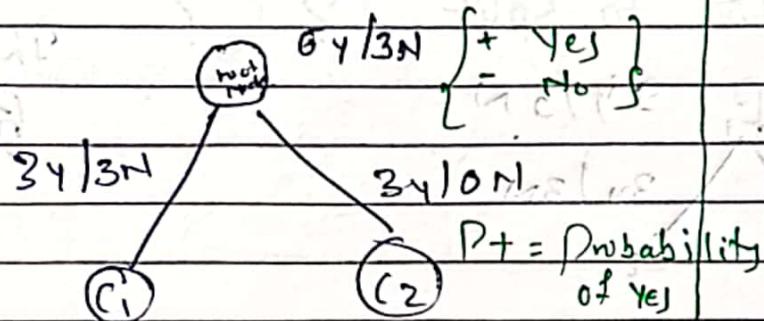
$\Rightarrow \Rightarrow$

1. Entropy

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

2. Gini Impurity

$$GI = 1 - \sum_{i=1}^n (P_i)^2$$

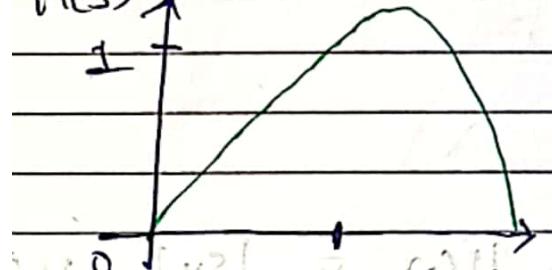


$$\text{Entropy} = H(S) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= -1 \log_2 1$$

$= 0 \rightarrow$  Pure split (100%).

$H(S)$



$$\left. \begin{aligned} P^+ &= 0.5 \\ P^- &= 0.5 \\ P &= 1 - Q \\ Q &= 1 - P \end{aligned} \right\}$$

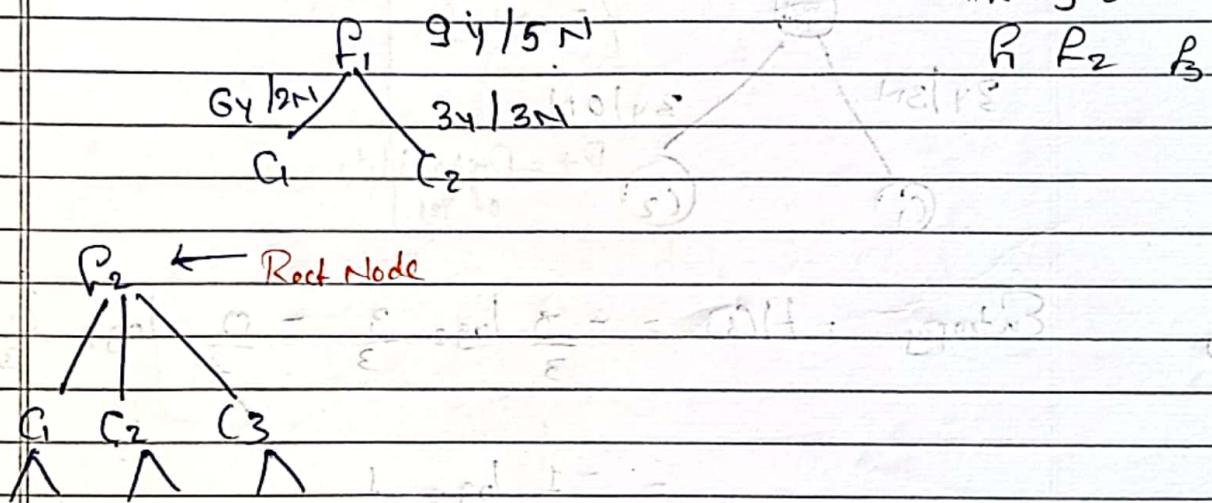
Pure split (100%)

$$\begin{aligned}
 \Rightarrow H(Q) &= -\frac{3}{6} - \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\
 &= -\frac{1}{2} - \log_2 \frac{3}{8} - \frac{1}{2} \log_2 \frac{1}{2} \\
 &= \underline{\underline{1}} \quad \rightarrow \text{Impure split } 50\% - 50\%
 \end{aligned}$$

$\Rightarrow$  Your always Entropy between 0 to 1

$\rightarrow$  ② Which feature to Split ? ? ?

In my data we

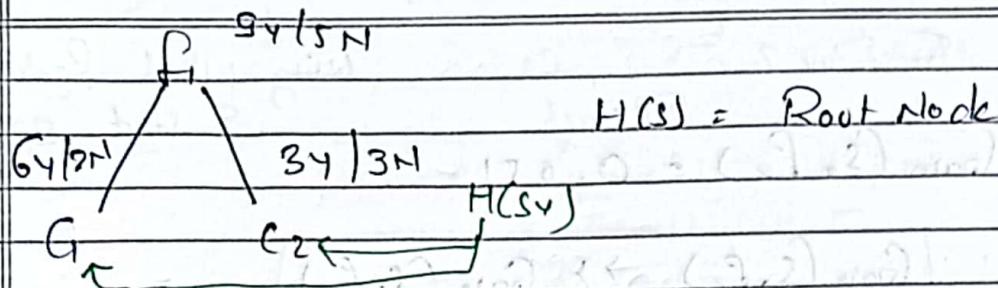


we basically we "Information Gain"

Formula :-

Gain (S)

$$\text{Gain (Sample ; P)} = H(S) - \sum_{V \in Values} \frac{|S_V|}{|S|} H(S_V)$$



$$H(S) = -P_1 \log_2 P_1 - P_2 \log_2 P_2$$

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \left( \frac{5}{14} \right)$$

$$\approx 0.94$$

For G<sub>1</sub>

$$H(S_{G_1}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$\boxed{H(G_1) = 0.81}$$

For C<sub>2</sub>

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

$$\boxed{H(C_2) = 1}$$

Finally,

$$\text{Grain}(S, F_1) = 0.94 - \left[ \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\text{Grain}(S, F_2) = \underline{\underline{0.48}}$$

For Feature 2 :-

$$\text{Gain}(S, P_2) = 0.651$$

Using which Feature should I start splitting first

$$\boxed{\text{Gain}(S, P_2) >> \text{Gain}(S, P_1)}$$

$P_2$  → to start the split

\* Gini Impurity :-

Formula :-

$$G.I. = 1 - \sum_{i=1}^n (P_i)^2$$

n = no of outputs

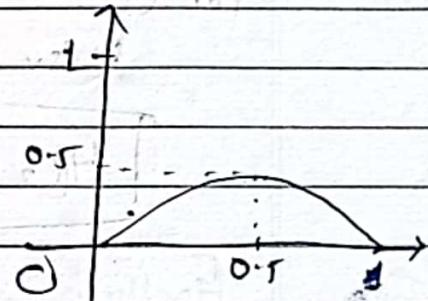
we have node 2 1/2

when Entropy = 1  
when Gini = 0.5

$$= 1 - [(P_+)^2 + (P_-)^2]$$

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

$$= 1 - \left[ \frac{1}{4} + \frac{1}{4} \right]$$



$$= \frac{1}{2}$$

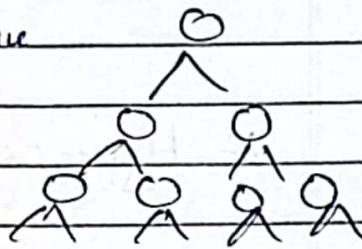
$$= 0.5$$

(Ans)

When we use Gini impurity (or) Entropy

→ More time where it will probably take  
for Execution purpose

→ Decision tree having a worst time complexity  
like 100 feature



Entropy → log Function  
Gini impurity → Simple maths

⇒ More execution time "Entropy"

When Small no of Feature then it will use  
Gini impurity is best

Fast

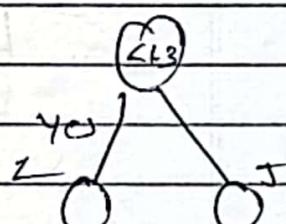
Gini >> Entropy

Let I Example Example

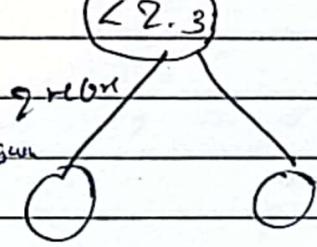
↪ Continuous

P	0.1P	$\Rightarrow f_1$
2.3	1.3	
1.3	2.3	
4	3	
5	4	
3	5	
7	7	
3		

Information gain



Leaf Node



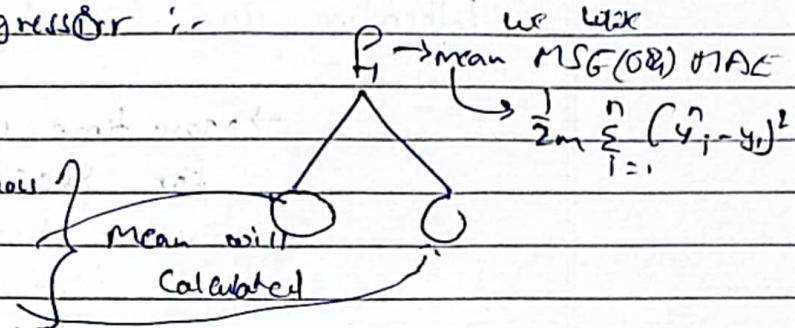
Best Information gain

Decision Tree regressor :-

and  $f_1, f_2$  will be  $O/P$ , and we will calculate

Continuous

Sum



\* Hyperparameters :-

Decision

Post Prunning

80%

74 / 2N

Pre-Prunning

(Stop)

Prunning decided by Hyperparameters

Max. depth, max. leaf.

→ GridsearchCV

What is the main problem of Decision Tree

Low Bias

High variance

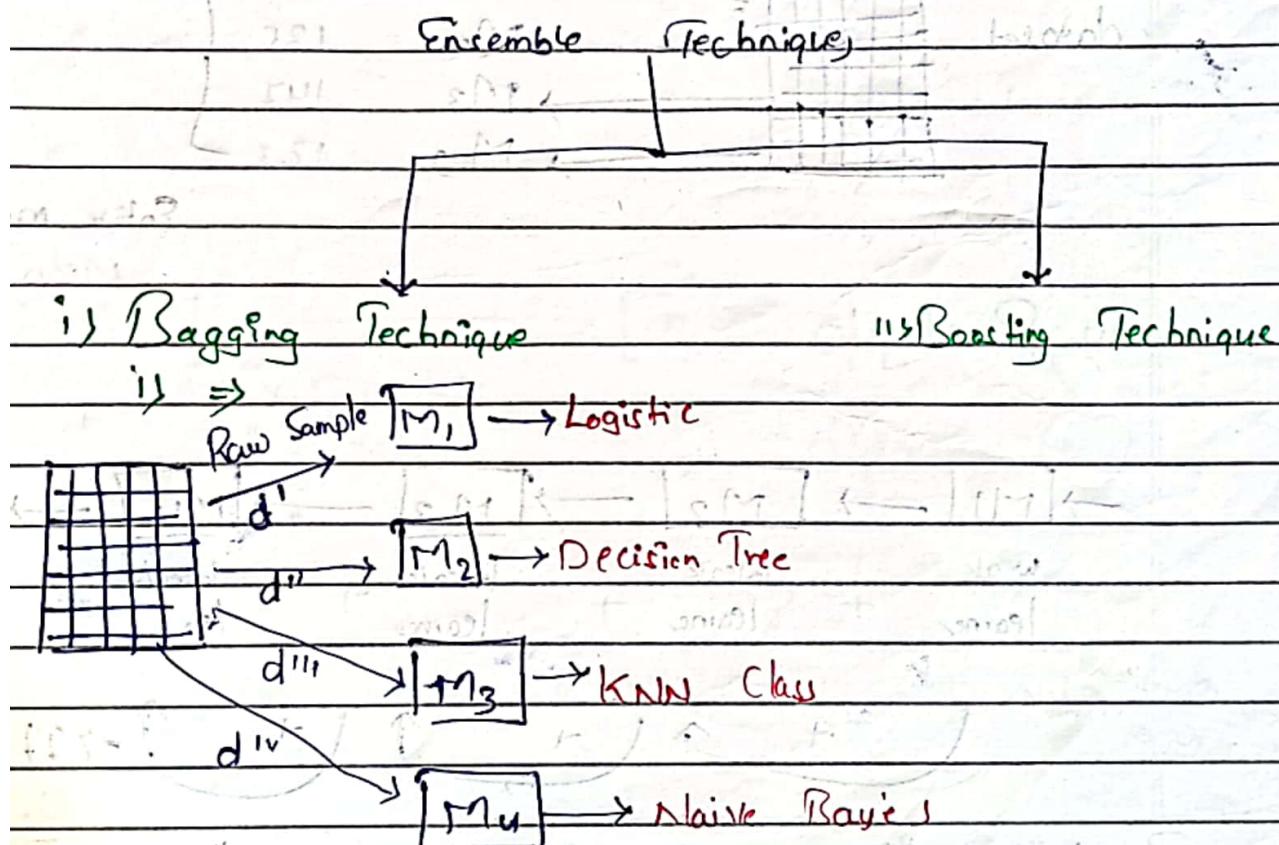
Overfitting ??

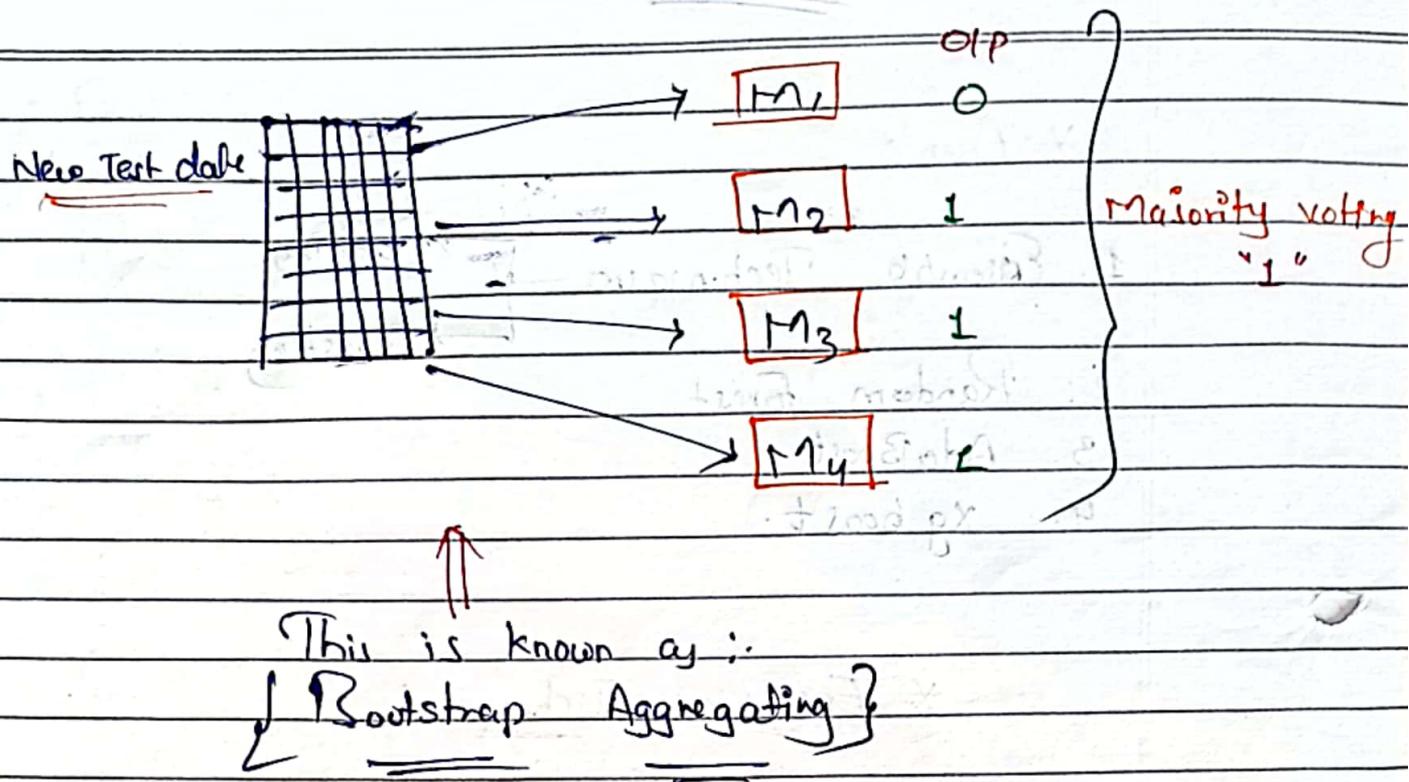
### \* Agenda

1. Ensemble Techniques
  - Bagging
  - Boosting
2. Random Forest
3. AdaBoost
4. Xgboost

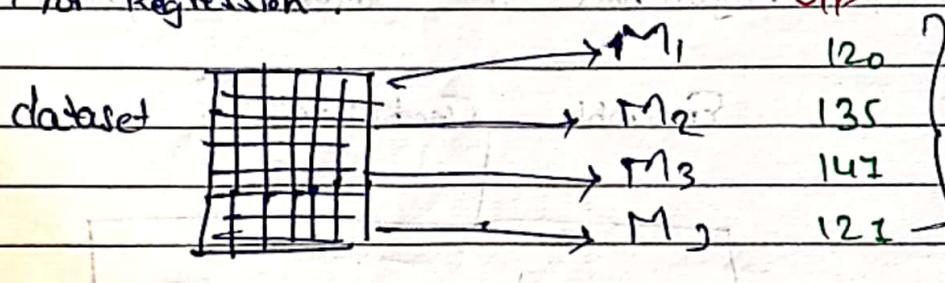
### \* Ensemble Techniques :-

Multiple Algorithms to Solve a problem ???

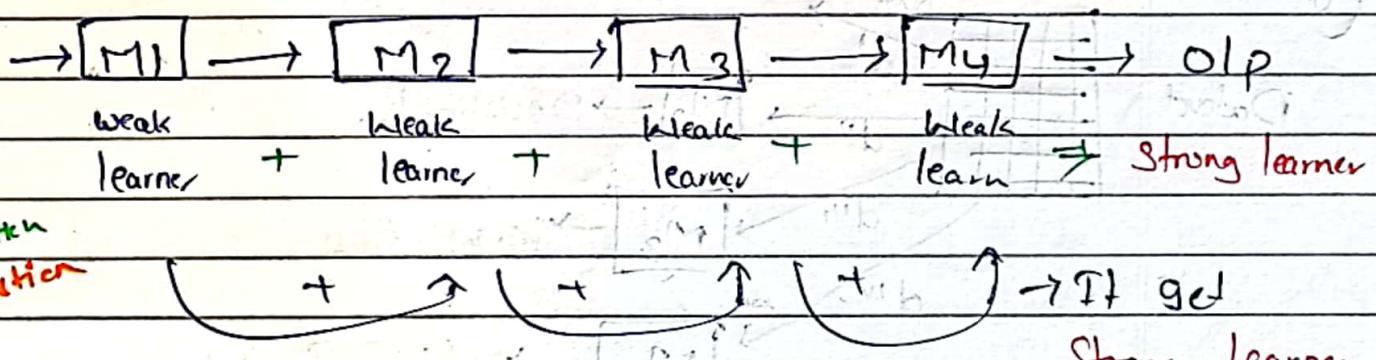




For Regression :-

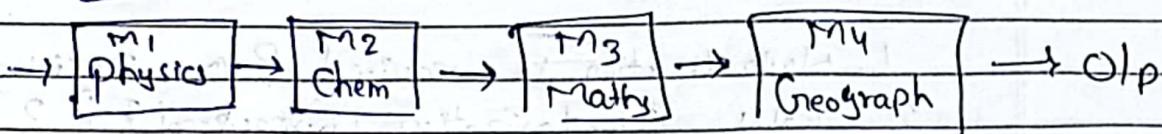


ii) Boosting :-



Boosting :- it is a Sequential set of all the models combined and these all models that i initialize are usually weak learners. When they are combined together they become a strong learner and

→ Example i -



⇒

Bagging

1. Random Forest classifier
2. Random Forest Regressor

Boosting

1. Ada boost
2. Gradient boost
3. Xg boost

⇒

• Random Forest classifier and Regressor :-

Random Forest  
dataset

here all models call DecisionTree

Row Sampling  
+  
Feature Sampling

[M<sub>1</sub>] → O/p → 0

[M<sub>2</sub>] → O/p → 1

[M<sub>3</sub>] → O/p → 0

[M<sub>4</sub>] → O/p → 0

majority  
→ O/p → 0

{ ⇒ for Regression to calculate mean of model output }

## Important point of in Random Forest :-

1. Normalization is required ? (or) Decision Tree  $\Rightarrow$  No
2. kNN (standardization) is required ? ?  $\Rightarrow$  Yes  
Euclidean distance  
Manhattan distance
3. Random Forest impacted by outliers ? No
4. kNN is impacted by outliers ? Yes.

Bagging  $\Rightarrow$  Random Forest

You can create custom bagging.

\* Boosting

is AdaBoost :- (Adaptive Boosting) :-

I have Dataset like

$F_1 F_2 F_3 F_4$  O/P Weight New update We define a weight  
--- --- --- --- Yes  $1/7$  0.9

No  $1/7$  0.05  $\rightarrow$  Information gain

$1/7$  0.05 Stumps Gini (or) Entropy

$1/7$  0.05  $P_1 P_2 P_3$

$1/7$  0.895 weak learner

$1/7$  0.05 weak learner

Total calculate Error =  $\frac{1}{7}$

② Performance of Stumps

Formula :-  $\frac{1}{2} \log_e \left( \frac{1 - \text{Total Error}}{\text{Total Error}} \right)$

$$\frac{1}{2} \log_e \left( \frac{1 - 1/7}{1/7} \right)$$

$$= 0.895$$

③ New Sample weight = <sup>update all weight & my dataset</sup>  
 Formula is weight  $\times e^{-Ps}$

Formula = weight  $\times e^{-Ps}$

( $Ps$  : Performance of stamp)

$$\text{New weight} = \frac{1}{7} \times e^{-0.895}$$

$$= 0.05$$

For,

$$\text{Incorrect records} = \text{weight} \times e^{Ps}$$

$$= \frac{1}{7} \times e^{0.895}$$

$$= 0.349$$

$$\therefore \left\{ P = \text{Approx } 2.71828 \right\} =$$

New weight

0.05 10.6us

0.05 10.6us

0.05 10.6us

0.05 0.349

0.05 10.6us

0.05 10.6us

0.05 10.6us

0.6us

It is 1 ??

Normalize weight.

$\frac{0.05}{0.07}$

0.07

0.07

$0.07 \frac{0.349}{0.07}$

0.07

0.07

0.07

0.07

Buckets

0 - 0.07

0.07 - 0.14

0.14 - 0.21

0.21 - 0.287

0.287 - 0.351

0.351

$k = 1$

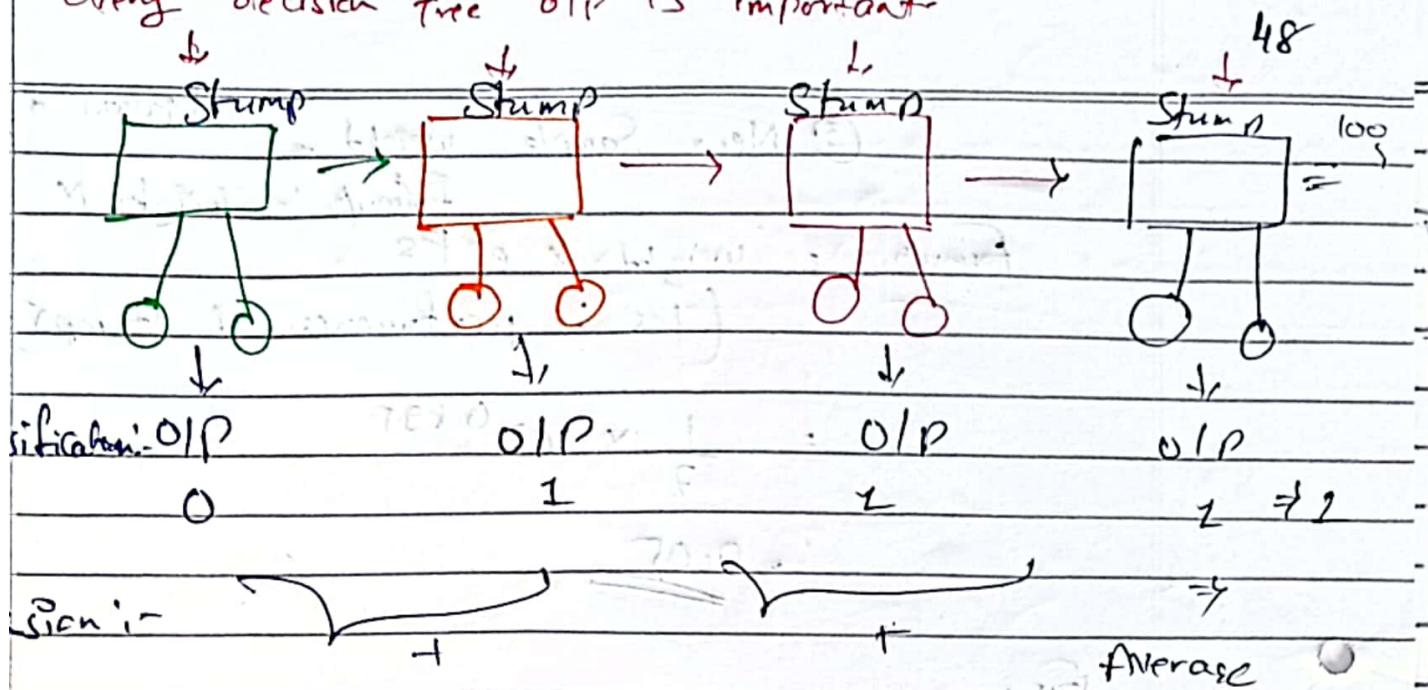
Bucket

Sum = 0.537

0.21

0.747

Every decision tree OIP is important



X Agendas :-

Un-supervised ML

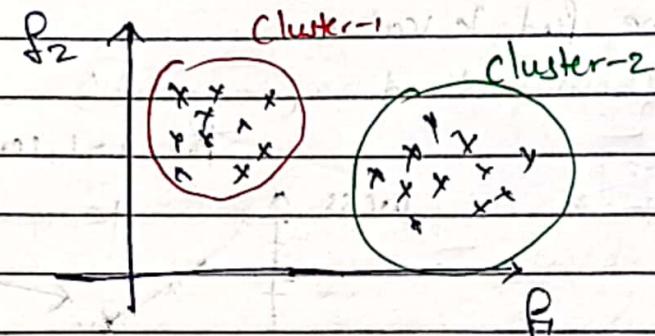
1. K Means clustering
2. Hierarchical Clustering
3. Silhouette Score
4. DBScan Clustering

Unsupervised ML  
→ We don't have specific O.P.

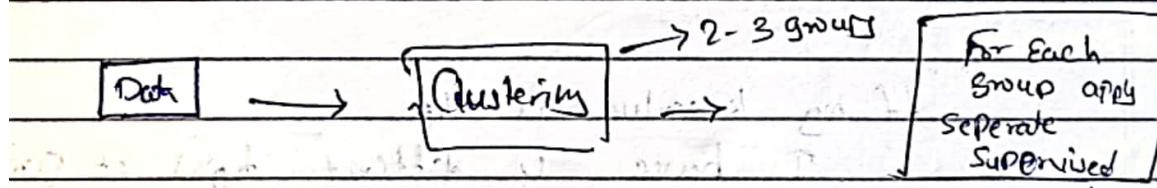


K-means clustering :-

Suppose we have datapoints

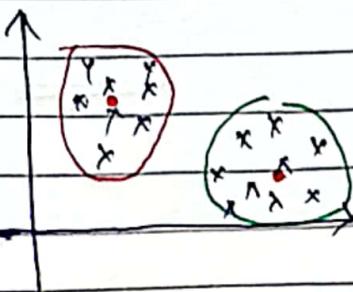


⇒ Custom Ensemble Techniques :-



K-Means

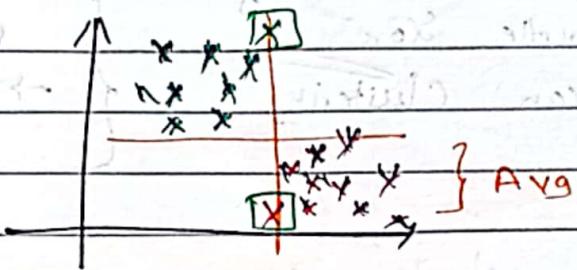
K = Centroid



Centroid

Step-1 :- Take try  $k$  values  $\rightarrow$  Suitable let  $K=2$

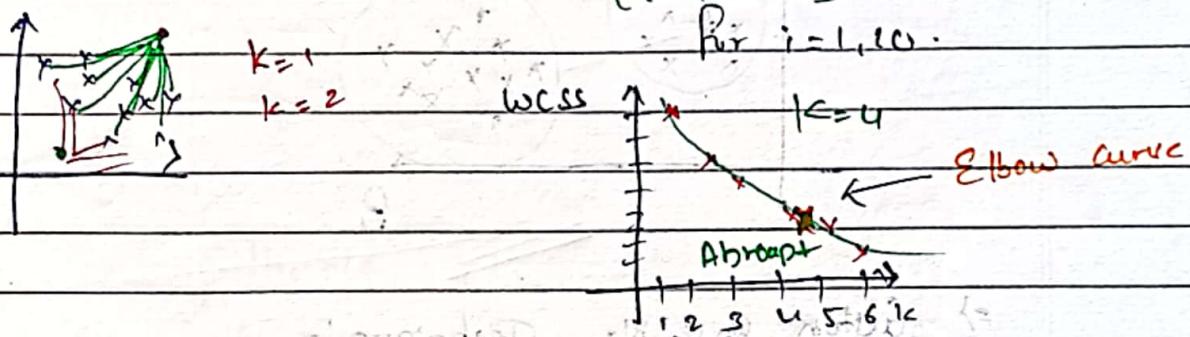
Step-2 :- Initialize  $k$  number of centroids



Step-3 :- Compute average to update centroid

$\Rightarrow$  How we find  $k$  value

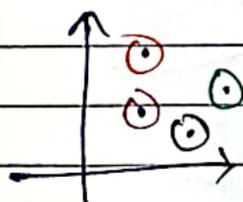
Elbow method ( $k$  value)



(WCSS :- within cluster sum of square)

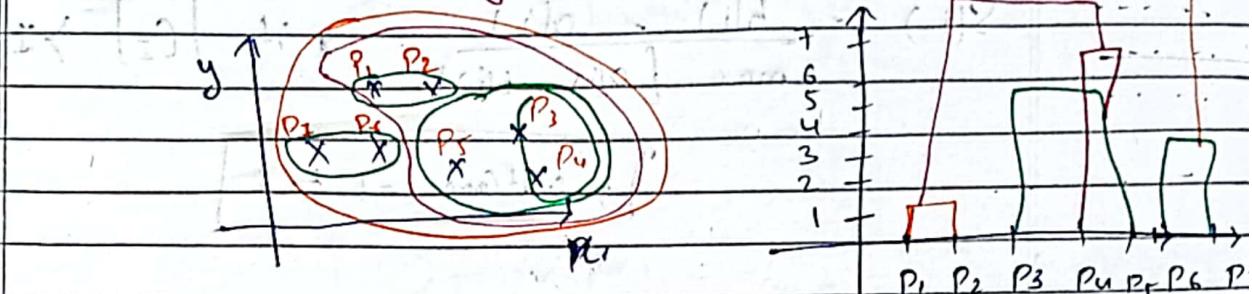
If my  $k$  value is 4

I have 4 different types of group  
and Every group have centroid



Dendrogram

## 2. Hierarchical clustering :-



{ you need to find the longest vertical line  
That has no horizontal tie past through it }

Q: Maximum Time is taken by K-means or  
Hierarchical Clustering ??

→ Hierarchical clustering will take more time

→ If dataset is small go with Hierarchical clustering

→ If dataset is large go with Clustering.

\* Validate clustering models :-

## Silhouette clustering

Calculation Average of Distance

$a(i)$

Cluster

Centroid

$b(i)$

Calculating

Good cluster i :-

$$a(i) \gg b(i)$$

Good model -  $b(i) \gg a(i)$  ✓✓

Formula:

$$S(i) := \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \text{ if } |C_i| > 1$$

$\Rightarrow \text{outCome} = 1 \text{ to } 1$

### DBScan clustering

Density-Based Spatial clustering of Applications with Noise (DBSCAN)

i. Min points

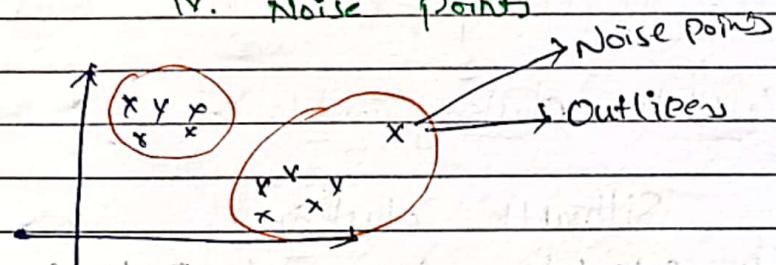
ii. Core point

iii. Border point

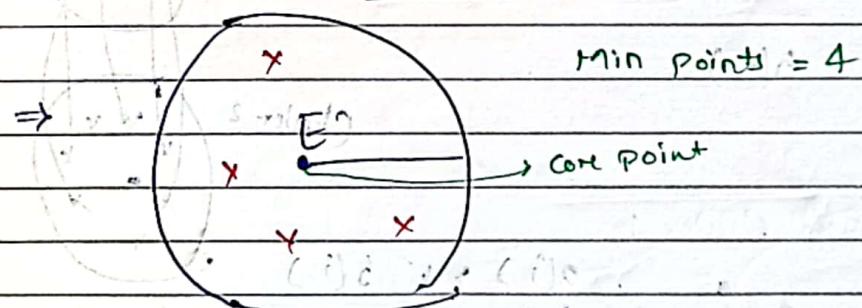
iv. Noise points

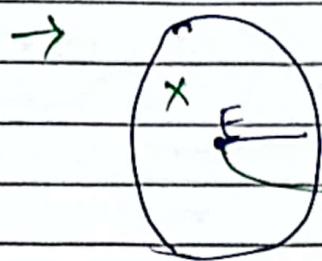
\* Epsilon

(radius of circle)

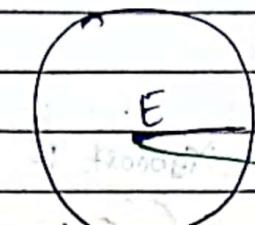


$\rightarrow$  this kind of algorithm skip the outliers we use DBScan





border point



## Noise point

18/01/25

## Day 7

54

### \* Agenda

1. Xgboost classifiers
2. Xgboost Regressor
3. SVM
4. SVR

Xgboost :- Full Form :-

1. Xgboost classifiers (- Extreme Gradient Boosting) :-

⇒ Dataset :-

Salary	Credit	Approval	Residra (Approval - Basemodel)
$\leq 50$	Bad	0	-0.5 (0 - 0.5)
$\leq 50$	Good	1	0.5 (1 - 0.5)
$\leq 50$	Good	1	0.5 (1 - 0.5)
$> 50$	Bad	0	-0.5 (0 - 0.5)
$> 50$	Good	1	0.5 (1 - 0.5)
$\geq 80$ k	Normal	1	0.5 (1 - 0.5)
$\leq 80$ k	N	0	-0.5 (0 - 0.5)

It is just a dummy :-

Base Model → Probability = 0.5

Step 1:- Create a binary Decision tree using the features

Step 2:- Calculate the similarity weight

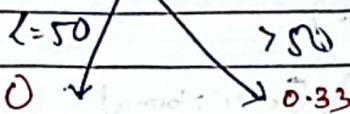
$$\text{Formula : } \frac{\sum (\text{Residra})^2}{\sum (\text{Probability} (1 - P))}$$

Step 3:- Information gain

vision

### Binary Decision [Salary]

Tree



$$[ -0.5, 0.5, 0.5, 0.5 ] \quad [ -0.5, 0.5, 0.5 ]$$

Now, we calculate similarity weight:-

$$\leq 50 :-$$

$$= \frac{(-0.5 + 0.5)(1-0.5)}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)}$$

$$= \underline{\underline{0}}$$

$\geq 50$ ,

$$\geq 50 :-$$

$$= \frac{(-0.5 + 0.5)(1-0.5)}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)}$$

$$= \frac{0.25}{0.75} = \frac{1}{3} = \underline{\underline{0.33}}$$

Now, by calculating Information gain

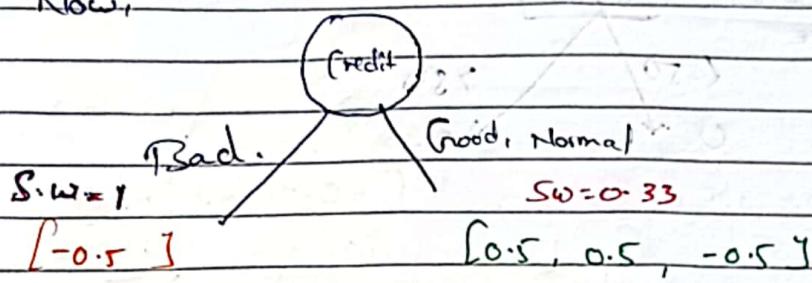
$$\text{Similarity weight} = \frac{0.25}{1.75} = \frac{1}{7} = \underline{\underline{0.142}}$$

Now Calculate Information gain

$$\text{Gained Information} = 0 + 0.33 - 0.14$$

$$= 0.19$$

Now,



Now,

Similarity weight is

$$\frac{0.25}{0.25} = 1$$

$$\frac{0.25}{0.25} = 1$$

Now,

$$\frac{0.25}{0.75} = \frac{0.33}{1}$$

$\Rightarrow$

$$1 + 0.33 - 0 \Rightarrow 1.33$$

(How we calculate the real probability :-

$$\log(P) - \log\left(\frac{P}{1-P}\right) = \log\left(\frac{0.5}{0.5}\right) = 0$$

$$r \cdot [0 + \lambda(1)]$$

$$[\lambda = 0.01]$$

Learning rate

$$O/P \rightarrow 1$$

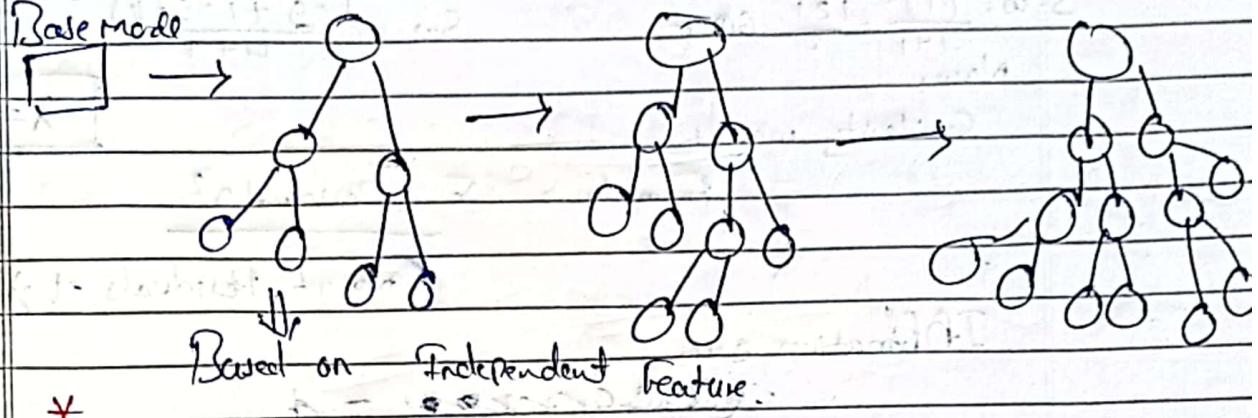
New Record

$$\left[ r_0 + \lambda_1(DT_1) + \lambda_2(DT_2) + \lambda_3(DT_3) + \lambda_4(DT_4) + \dots + \lambda_n(DT_n) \right]$$

Xgboost  $\rightarrow$  Black box model.

$\lambda \Rightarrow$  Cross validation

$$\text{Base model} + \lambda_1(DT_1) + \lambda_2(DT_2) + \dots + \lambda_n(DT_n)$$



Xgboost Regressor :-

Exp	Graph	Salary	Residual
2	Yes	40k	-11k
2.5	Yes	42k	-10k
3	No	52k	1k
4	No	60k	9k
4.5	Yes	62k	8.11k

Now, we create a Base model:-

Base model  $\rightarrow$  Avg 51k

$$[-11, -9, 1, 9, 11] \rightarrow S_{WI} = \frac{(-11 + 9 + 1 + 9 + 11)^2}{5+1} = \frac{1}{6}$$

Exp

$\leq 2$

$> 2$

$60 \cdot 5$

$S_{WI} = 28.8$

$$\therefore [-11] \quad [-9, 1, 9, 11] \\ w = \frac{(-11)^2}{1+1} = \frac{121}{2} = 60.5$$

$$S_{WI} = \frac{(-9 + 1 + 9 + 11)^2}{4+1} = \frac{144}{5} = 28.8$$

Now,

$\lambda = 0$

Similarity weight

$$\text{formula} = \frac{\sum (\text{Residual})^2}{\text{No. of Residuals}}$$

$\Rightarrow 60.5 + 28.8 = \underline{\underline{1/6}}$

$$= \underline{\underline{89.13}}$$

$$= 60.5 + 28.8 - \underline{\underline{1/6}}$$

$$= \underline{\underline{89.13}}$$

Now,

Exp

$\geq 2.5$

$2.5$

$S_{WI} = 110.25$

$$= 133.33 \quad [-11, -9] \quad [1, 9, 11]$$

$$S_{WI} = \frac{(-11 + 9)^2}{2+1}$$

$$= \frac{(-20)^2}{3} = \underline{\underline{133.33}}$$

$$S_{WI} = \frac{(1 + 9 + 11)^2}{1+3}$$

$$\frac{144}{4}$$

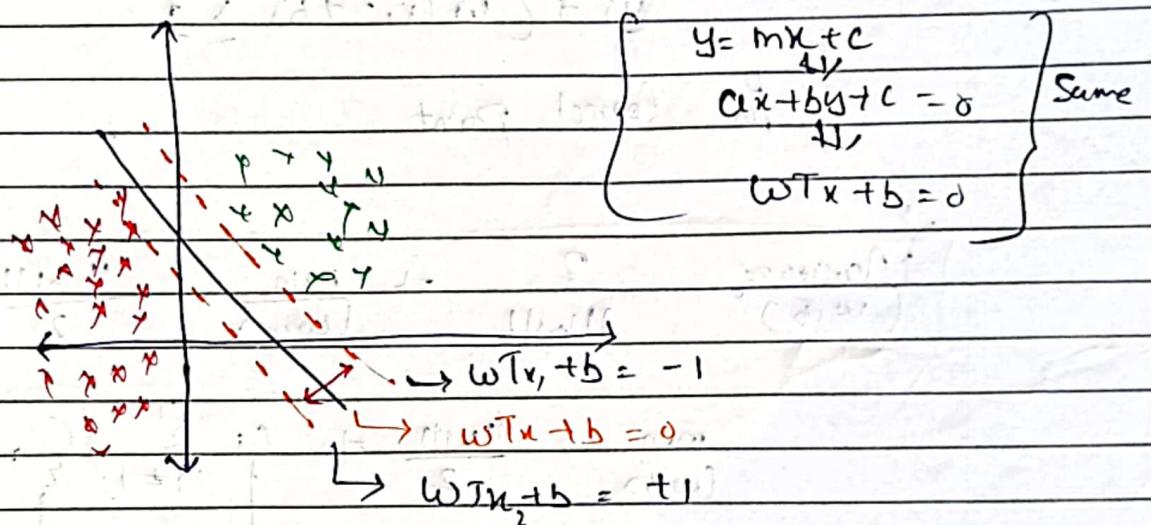
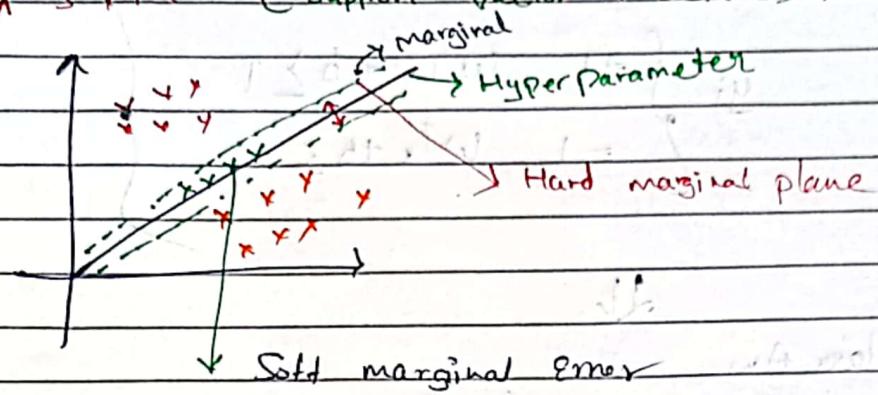
$$\frac{144}{4} = \underline{\underline{110.25}}$$

Now, here,

Now, here

$$51 + \alpha_1(-1) - \alpha_2 + \alpha_2(\alpha_1 + \alpha_2) \quad \left| \begin{array}{l} 51 + \alpha_2(7) \end{array} \right.$$

\* SVM (Support vector machine) :-



$$w^T x_1 + b = 1$$

$$w^T x_2 + b = -1$$

(Vector + magnitude)

$$w^T(x_1 - x_2) = 2$$

$$\frac{|w|}{|w|}$$

60

$$\text{Maximize}_{(w,b)} \frac{2}{\|w\|} \quad \downarrow \text{Maximize loss function} \quad \rightarrow \text{marginal plane}$$

Such that,

$$y_i \left\{ \begin{array}{l} +1 \quad w^T x + b \geq 1 \\ -1 \quad w^T x + b \leq -1 \end{array} \right.$$

$$\text{Margin} \quad y_i (w^T x_i + b) \geq 1$$

for Correct Point

$$\text{Maximize}_{(w,b)} \frac{2}{\|w\|} \Rightarrow \min_{(w,b)} \frac{\|w\|^2}{2}$$

$$\min_{(w,b)} \frac{\|w\|^2}{2} + C_i \sum_{i=1}^n \{ \}$$

(Summation of the distance of the wrong point)

\* Sym Kernel's

How many errors

we can have ~

