

"Logistic Regression is a classification Problem" Date: \_\_\_\_\_

## LOGISTIC REGRESSION:

- In classification, the output 'y' is discrete value
- In logistic regression, to develop an algorithm to determine what class a new input should fall into.

Example: Email: Spam/not Spam

PROBLEM WE DISCUSS → Binary Classification (0 or 1)

- How to develop the classification algorithm?
- If we use linear regression there exists a problem that if the number is greater than 1 or less than 0 so, in which category we will classify it.
- Linear regression donot work on mixed data as well as Linear regression could not solve boundary problem.
- Hypothesis can give value larger than 1 or less than 0

## LOGISTIC REGRESSION (CLASSIFICATION)

### HYPOTHESES REPRESENTATION:

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

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

so,

$$g(z) = g(\underline{\theta^T x})$$

Therefore,

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

( $\because g(z) = \text{sigmoid}$ )

$\therefore z$  is a real number

Where  $z = \underline{\theta^T x}$

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

$$\begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$\begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$\theta_0 + \theta_1 x = \underline{\theta^T x}$$

$$\text{Hypotheses} = h(x) = \frac{1}{1 + e^{-\alpha x}}$$

This is the Sigmoid function or the logistic function.

Question:

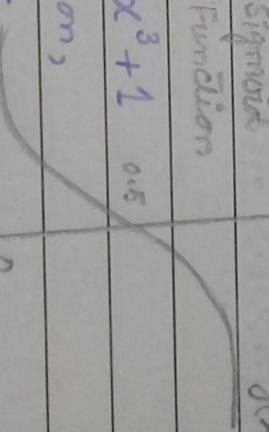
- What value should be of  $\alpha^T x$ ?
- It is not necessary that the value of  $\alpha^T x$  is always linear ( $\alpha_0 + \alpha_1 x$ ), it can be:  

$$\alpha^T x = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots$$
- or  

$$\alpha^T x = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2^2 + \dots$$
- What decides the value of  $x$

The value of  $x$  will be decided by regularization

### SIGMOID FUNCTION:

- Crosses 0.5 at the origin, then flattens out
  - Asymptotes between 0 and 1
  - Asymptotes means: They go towards or give Nearest value 0 or 1, but not exactly reach to 0 or 1 (decimal value)
- |                  |            |
|------------------|------------|
| Sigmoid Function | $y = g(x)$ |
|------------------|------------|
- 

- $g(\alpha_0, \alpha_1) \rightarrow$  similar to  $f(x) = x^2 + x^3 + 1$  at 0.5
- In  $f(x)$   $x$  is a part of a function,
- similarly, in  $g(\alpha_0, \alpha_1)$   $\alpha$  is a part of a function.

- KNN for classification & Regression
- Understanding How it works
- Example
- Python code

K-Nearest Neighbour used for both Linear & Classification

It is the lazy algorithm:

- because when you work on linear regression

→ When you develop model there exists a plan or best line or a decision boundary behind the model

→ In KNN the interpretation of model of training data load into memory & when you give test data then it reads from initial, hitting training data again & again, we don't make model in KNN

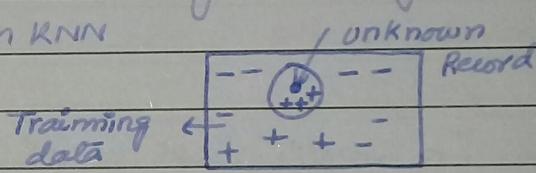
Nearest Neighbour Classifier

Requires three things

→ The set of stored records

→ Distance Metric to compute distance b/w records.

→ The value of K, the number of nearest neighbours to retrieve



- لـ K Method کو چیزی کہ اسے distance سے ایک نزدیکی کرنا  
- create distance میں مبتداً

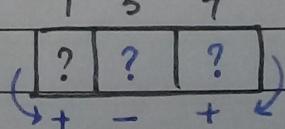
(K) → nearest کا ذمہ دینا  
- مابینہ میتوں

1 test sy n distance agy

نیکے دیکھا قریب ترین

IF K=3 the we pick<sup>3</sup> nearest distances

small



Normally K is odd

# Sign Similarity } Classification Majority Voting

Date: \_\_\_\_\_

- We are using training data again & again for that's why it is called "Lazy algorithm."
- For linear regression, if we apply KNN so, what happened?
- Calculate distance
- Store distance
- 5 Nearest Neighbour apply
- Take Mean

4	8	11	19	51
?      ?      ?      ?      ?				

Value      sum of all values / total  
mean value

distance Metric → kons a tool  
for similarity

Sno	X <sub>1</sub>	X <sub>2</sub>	Y
1	2	5	3.9
2	6	3	4.1
3	7	4	2.7
4	1	1	1

10      8      3      ?

→ cosine

→ Kaisa lo k tie na ho

→ We will take mean if there is a tie

→ minimum distance with majority voting

→ In multiclass problem so, how to decide target class?

→ majority voting

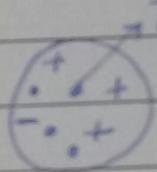
→ The minimum mean result

→ It is used for Partitioning

Measuring tool = Euclidean distance

The formula of KNN by Euclidean Distance

$$Dist(C_1, C_2) = \sqrt{\sum_{i=1}^n (attr_i(C_1) - attr_i(C_2))^2}$$



Multi-class Problem  
We have 2 attributes

Date: \_\_\_\_\_

	Sweet	Crunch	Food Type
Grape	8	5	Fruit
Greenbean	3	7	Vegetable
Nuts	3	6	Protein
Orange	7	3	Fruit
Tomato	6	4	?

$$\text{Maximum } K = 4$$

$$\text{Minimum } K = 1$$

$$\text{Dist}^{(1)}(C_T, C_1) = \sqrt{(6-8)^2 + (4-5)^2} \\ = \sqrt{5}$$

$$\text{Dist}^{(1)}(C_T, C_1) = 2 \cdot 2$$

$$\text{Dist}^{(2)}(C_T, C_2) = \sqrt{(6-3)^2 + (4-7)^2} \\ = \sqrt{18}$$

$$\text{Dist}^{(2)}(C_T, C_2) = 4 \cdot 2$$

$$\text{Dist}^{(3)}(C_T, C_3) = \sqrt{(6-3)^2 + (4-6)^2} \\ = \sqrt{13}$$

$$\text{Dist}^{(3)}(C_T, C_3) = 3 \cdot 6$$

$$\text{Dist}^{(4)}(C_T, C_4) = \sqrt{(6-7)^2 + (4-3)^2} \\ = \sqrt{2}$$

$$\text{Dist}^{(4)}(C_T, C_4) = 1 \cdot 4$$

$$K = 1$$

$$\min_1 (\text{Dist}^{(1)}, \text{Dist}^{(2)}, \text{Dist}^{(3)}, \text{Dist}^{(4)})$$

Now,  $\text{Dist}^{(4)} \rightarrow \text{Fruit}$

$\downarrow$   
min

Prediction for Tomato  
Fruit

FOR K=2

Dist<sup>(1)</sup>  
 $\downarrow_F$  , Dist<sup>(3)</sup>  
 $\downarrow_P$

K = 4

Dist<sup>(1)</sup>    Dist<sup>(2)</sup>    Dist<sup>(3)</sup>    Dist<sup>(4)</sup>  
|            |            |            |  
F            P            P            F

→ We will took out mean & then, we classify  
Outcome → will lies where  
load data  
in

## LINEAR REGRESSION (IMPORTANT POINTS)

- Supervised Learning algorithm + linear regression is not uncertain
- Works on Continuous values
- draws a linear line

## COST FUNCTION FORMULA & GRAPH<sub>2</sub>

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

Cost function is used to find  $\theta$  only

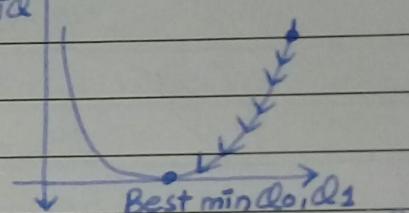
## GRADIENT DESCENT FORMULA & GRAPH<sub>2</sub>

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

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

General formula of Gradient Descent:

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



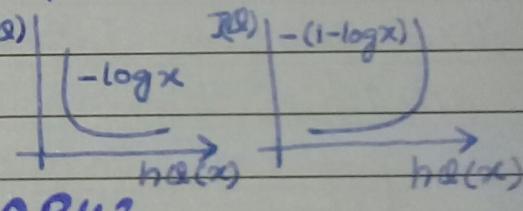
Gradient Descent is used to find best minimum  $\theta$

## LOGISTIC REGRESSION (IMPORTANT POINTS)

→ Supervised learning algorithm

→ Works on discrete Values

→ draws a boundary line b/w classes.



## COST FUNCTION FORMULA & GRAPH<sub>2</sub>

$$J(\theta) = -\frac{1}{m} \left[ \underbrace{\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)})}_{\text{FOR 1}} + \underbrace{(1-y^{(i)}) \log (1-h_\theta(x^{(i)}))}_{\text{FOR 0}} \right]$$

Cost function is used to find  $\theta$  only

## GRADIENT DESCENT FORMULA & GRAPH<sub>2</sub>

$$\theta_j = \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

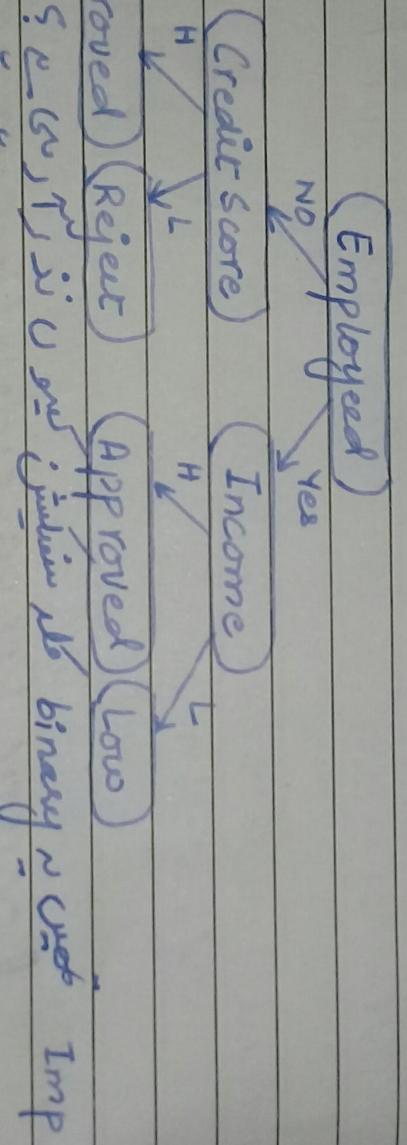
## DECISION TREE

It is used for classification & regression.

→ It is mostly used for classification.

→ Algorithm used for decision tree ID3 Heuristic algorithm.

- A decision tree is a tree structured classifier
- It has two types of Nodes
  - decision nodes
  - leaf nodes



### Decision Tree Introduction

attribute → mean nodes

The popular attribute selection measure

Information gain (Pink box / Blue box)

Gini index (Python & R uses Gini method)

Hum alone ko use kr skte hain

Has attribute may kloni information Hai

Play is Categorical so, we solved it with Information gain

→ This may sub sy riyada information Hogi usko hum Parent Node banyay

→ Or gay nikalnay k liyay num entropy use kryngay i.e.  $H = -P(+)\log_2 P(+)$

→ Or recursively

Tree binary Tree hora jawai nahi Hai

What is Pure?

→ Pure against attributes

- e.g. 0 or 1

→ Impure against attributes mix tree

attributes khtm hajnor pure na ho to gay

jo eror ala

→ ID3 is used for making Tree

- What is the importance of node?

Decision Tree:

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

Entropy Probability Yes  
Representation No

→ S... Subset of training examples

→  $P(+)$  or  $P(-)$  ... % of positive or negative examples in S

Cer

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

$$= \frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6}$$

Date: \_\_\_\_\_

How we gain information Gain  
→ How many items in pure sets

→ main entropy

$$\text{Gain}(S, A) = H(S) - \sum_{\text{VE values}(A)} \frac{|S_{V_i}|}{|S|} H(S_{V_i})$$

→ V: possible values of A

→ S: set of examples (x)

→ SV: Sub set where  $X_A = V$   
9 Yes / 5 No

$$-\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2$$

Wind

$$\frac{5}{14}$$

$$H(S) = 0.94$$

Weak

Strong

Gain(S, Wind)

$$= H(S) - \frac{8}{14} H(S_{\text{weak}}) - \frac{6}{14} H(S_{\text{weak}})$$

$$= 0.94 - \frac{8}{14} * 0.81 - \frac{6}{14} * 1.0$$

$$= 0.049$$

First we find main Information Gain

Date: \_\_\_\_\_

$$-\frac{19}{4} \log_2 \frac{9}{14}$$

Information gain & Entropy

Example To find Parent node among  
Step # 01: Select parent node  
:- calculate IG

Step # 01: To find parent node among the provided nodes

i - to calculate IG for all nodes

ii - to calculate entropy

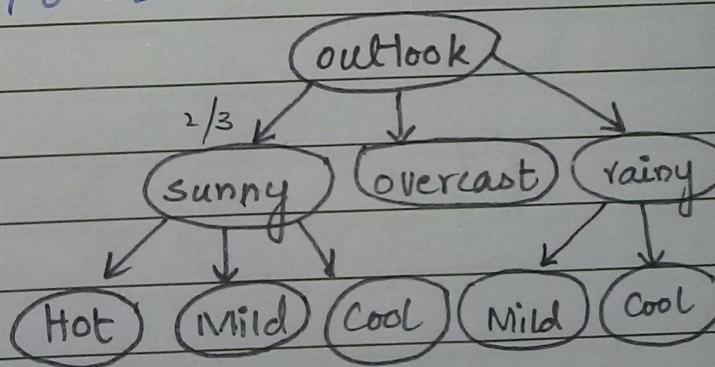
outlook or split

$$H(S) = (3) = - \left[ \frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} \right]$$
$$= 0.409 + 0.530$$

$$9/14 = 14$$

$$HS(3) = 0.940$$

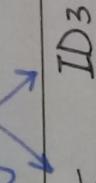
Imp



## DECISION TREES

→ A "Decision tree" is a tree which where each node represents a "Feature (Attribute)", each link (Branch) represents a "Decision (Rule)" and each leaf represents an "Outcome".

Algorithms



Gini Index

Entropy function

Information Gain

Sno	Outlook	Temperature	Humidity	Windy	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rainy	Mild	High	Weak	Yes
5	Rainy	Cool	Normal	Weak	Yes
6	Rainy	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rainy	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rainy	Mild	High	Strong	No

→ If we have ~~more~~ number of positive & negative than,  
Entropy = 0

→ ID we have only Positive Example & Negative Example  
Entropy = 0

Play tennis is a class attribute

→ Step 1: Create a Root Node?

• How to choose the attribute as the root Node?

The attribute that best classifies the training data,  
use this attribute as the root of the Tree.

→ Calculate **ENTROPY** (Amount of uncertainty in dataset)

$$\text{Entropy} = -\frac{P}{P+n} \log_2 \left( \frac{P}{P+n} \right) - \frac{n}{P+n} \log_2 \left( \frac{n}{P+n} \right)$$

→ Calculate **AVERAGE INFORMATION**:

$$I(\text{Attribute}) = \sum \frac{P_i + n_i}{P+n} \text{Entropy}(A)$$

→ Calculate **INFORMATION GAIN**: Difference in Entropy before &  
after splitting dataset on Attribute (A)

$$\text{Gain} = \text{Entropy}(S) - I(\text{Attribute})$$

ID3 Algorithm Steps:

- 1- Compute Entropy for data set Entropy (S)
- 2- For every attribute or Feature:
  - 1- Calculate entropy for all other values Entropy(A)
  - 2- Take average Information Entropy for the current attribute
  - 3- Calculate Gain for the current attribute
  - 3- Pick the Highest gain attribute
  - 4- Repeat until we get the desired tree.

Count Play Tennis Positive & Negative

Play Tennis	Yes	No	Total = 14
	9	5	

$$\text{Entropy} = -\frac{P}{P+n} \log_2 \left( \frac{P}{P+n} \right) - \frac{n}{P+n} \log_2 \left( \frac{n}{P+n} \right)$$

$$\text{Entropy} = -\frac{9}{9+5} \log_2 \left( \frac{9}{9+5} \right) - \left( \frac{5}{9+5} \right) \log_2 \left( \frac{5}{9+5} \right)$$

$$\text{Entropy} = -\frac{9}{14} \log_2 \left( \frac{9}{14} \right) - \frac{5}{14} \log_2 \left( \frac{5}{14} \right) = 0.940$$

→ For each attribute (let say Outlook)

- Calculate Entropy for each values: 'Sunny', 'Rainy', 'Overcast'

Outlook	Yes	No	Total	Outlook
* Rainy	3	2	5	
* Overcast	4	0	4	
* Sunny	2	3	5	

$$\text{Entropy of Rainy} = -\frac{3}{3+2} \log_2 \left( \frac{3}{3+2} \right) - \frac{2}{3+2} \log_2 \left( \frac{2}{3+2} \right)$$

$$\text{Entropy of Rainy} = 0.971$$

$$\text{Entropy of Sunny} = -\frac{2}{2+3} \log_2 \left( \frac{2}{2+3} \right) - \left( \frac{3}{2+3} \right) \log_2 \left( \frac{3}{2+3} \right)$$

$$\text{Entropy of Sunny} = 0.971$$

$$\text{Entropy of Overcast} = -\frac{4}{4+0} \log_2 \left( \frac{4}{4+0} \right) - \left( \frac{0}{4+0} \right) \log_2 \left( \frac{0}{4+0} \right)$$

$$\text{Entropy of Overcast} = 0$$

- Calculate Average Information Entropy:

$I(\text{Outlook}) = \frac{P_{\text{sunny}} + n_{\text{sunny}}}{P_{\text{th}}} \text{ Entropy} (\text{outlook} = \text{sunny}) +$

$\frac{P_{\text{rainy}} + n_{\text{rainy}}}{P_{\text{th}}} \text{ Entropy} (\text{outlook} = \text{rainy}) +$

$\frac{P_{\text{overcast}} + n_{\text{overcast}}}{P_{\text{th}}} \text{ Entropy} (\text{outlook} = \text{overcast})$

$$I(\text{Outlook}) = \frac{2+3}{9+5} (0.971) + \frac{3+2}{9+5} (0.971) + \frac{4+0}{9+5} (0)$$

$$I(\text{outlook}) = 5.3944 + 0.3467 + 0$$

$$I(\text{outlook}) = 0.693$$

- Calculate Gain: attribute is outlook

$$\text{Gain} = \text{Entropy}(S) - I(\text{Attribute})$$

$$\text{Gain} = 0.940 - 0.693$$

$$\text{Gain} = 0.247$$

Sunny	Yes	No	Total	Entropy
* Hot	0	2	2	
* Mild	1	1	2	
* Cool	1	0	1	

$$\text{Entropy of sunny \& Hot} = -\frac{0}{0+2} \log_2 \left( \frac{0}{0+2} \right) - \left( \frac{2}{0+2} \right) \log_2 \left( \frac{2}{0+2} \right)$$

$$\text{Entropy of sunny \& Hot} = 0$$

$$\text{Entropy of sunny \& Mild} = -\frac{1}{1+1} \log_2 \left( \frac{1}{1+1} \right) - \frac{1}{1+1} \log_2 \left( \frac{1}{1+1} \right)$$

$$\text{Entropy of sunny \& Mild} = 1$$

$$\text{Entropy of sunny \& Cool} = -\frac{1}{1+0} \log_2 \left( \frac{1}{1+0} \right) - \frac{0}{1+0} \log_2 \left( \frac{0}{1+0} \right)$$

$$\text{Entropy of sunny} = 0$$

• Calculate Average Information Entropy:

$$I(Sunny) = \frac{P_{Hot} + n_{Hot}}{P+n} \text{ Entropy (sunny=Hot)} +$$

$$\frac{P_{Mild} + n_{Mild}}{P+n} \text{ Entropy (sunny=Mild)} + \frac{P_{Cool} + n_{Cool}}{P+n} \text{ Entropy (sunny=Cool)}$$

$$I(\text{sunny}) = \frac{0+2}{5}(0) + \frac{1+1}{5}(1) + \frac{1+0}{5}(0)$$

$$I(\text{sunny}) = 0.4$$

• Calculate Gain: attribute is Sunny

$$\text{Gain} = \text{Entropy}(S) - I(\text{sunny})$$

$$\text{Gain} = 0.940 - 0.4$$

$$\text{Gain} = 0.54$$

Sunny and Humidity	Yes	No	Total	Entropy
High	0	3	3	
Normal	2	0	2	

$$\text{Entropy of sunny \& Humidity} = -\frac{P}{P+n} \log_2 \frac{P}{P+n} - \frac{n}{P+n} \log_2 \frac{n}{P+n}$$

= (High)

$$\text{Entropy of } S \& H = \text{High} = -\frac{0}{0+3} \log_2 \frac{0}{0+3} - \frac{3}{0+3} \log_2 \frac{3}{0+3}$$

$$\text{Entropy of } S \& H = \text{High} = 0$$

$$\text{Entropy of } S \& H = \text{Normal} = -\frac{2}{2+0} \log_2 \frac{2}{2+0} - \frac{0}{2+0} \log_2 \frac{0}{2+0}$$

$$\text{Entropy of } S \& H = \text{Normal} = 0$$

Calculate Average Information Entropy:

$$I(\text{Humidity}) = \frac{P_{\text{Humidity}} + n_{\text{Humidity}}}{P+n} \text{ Entropy (Humidity=High)} +$$

$$\frac{P_{\text{normal}} + n_{\text{normal}}}{P+n} \text{ Entropy (Humidity=Normal)}$$

$$I(\text{Humidity}) = \frac{0+3}{5} (0) + \frac{2+0}{5} (0)$$

$$I(\text{Humidity}) = 0$$

• Calculate Gain: attribute is Humidity

$$\text{Gain} = \text{Entropy}(S) - I(\text{Humidity})$$

$$\text{Gain} = 0.940 - 0$$

$$\text{Gain} = 0.94$$

Sunny & Windy	Yes	No	Total Entropy
Weak	1	2	3
Strong	1	1	2
Entropy of Sunny & Windy	$-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$		
= Weak	$\frac{1}{3} P_{\text{W}} \log_2 \frac{1}{3} + \frac{2}{3} P_{\text{W}} \log_2 \frac{2}{3}$		
Entropy of Sunny & Windy	$-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$		
= Weak	$\frac{1}{3} P_{\text{W}} \log_2 \frac{1}{3} + \frac{2}{3} P_{\text{W}} \log_2 \frac{2}{3}$		
Entropy of Sunny & Windy	$0.9182$		
= Weak			
Entropy of Sunny & Windy	$-1 \cdot \log_2 \frac{1}{2} - 2 \cdot \log_2 \frac{1}{2}$		
= Strong	$\frac{1}{2} P_{\text{W}} \log_2 \frac{1}{2} + \frac{1}{2} P_{\text{W}} \log_2 \frac{1}{2}$		
Entropy of sunny & Windy	$1$		
= Strong			

Calculate Average Information Entropy:

$$I(\text{Windy}) = P(\text{Weak}) + n \underline{\text{Weak Entropy}} (\text{Windy} = \text{Weak}) + P_{\text{W}} n$$

$$P(\text{strong}) + n \underline{\text{strong Entropy}} (\text{Windy} = \text{strong})$$

$$I(\text{Windy}) = \frac{1+2}{5} (0.9182) + \frac{1+1}{5} (1)$$

$$I(\text{Windy}) = 0.9509$$

$$I(\text{Windy}) = 0.55092 + 0.4$$

$\text{error}_{\text{train}}(h) < \text{error}_{\text{train}}(h')$

$\text{error}_{\text{test}}(h) > \text{error}_{\text{test}}(h')$

Date: \_\_\_\_\_

Calculate Information gain for Attribute (A)

$$H(S) = 1$$

$$\text{var } A >= 5$$

$$H(A >= 5) = -\frac{5}{5+7} \log_2 \frac{5}{12} - \frac{7}{12} \log_2 \frac{7}{12}$$

$$H(A >= 5) = 0.9799$$

$$H(A < 5) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}$$

$$H(A < 5) = 0.81128$$

$$IG_I(A) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)$$

$$IG_I(A) = 1 - \left[ \frac{12}{16} \times 0.9799 + \frac{4}{16} \times 0.81128 \right]$$

$$IG_I(A) = 0.0622$$

Information gain of D

$$H(D >= 1.4) = -\frac{0}{4} \log_2 \frac{0}{4} - \frac{1}{7} \log_2 \frac{1}{7}$$

$$H(D >= 1.4) = 0$$

$$H(D < 1.4) = -\frac{8}{9} \log_2 \frac{8}{9} - \frac{1}{9} \log_2 \frac{1}{9}$$

$$H(D < 1.4) = 0.5032$$

$$IG_I(D) = 1 - \left[ \frac{7}{16} \times 0 + \frac{9}{16} \times (0.5032) \right]$$

$$IG_I(D) = 1 - (0 + 0.2830)$$

$$IG_I(D) = 0.717$$

$$> = 4 \cdot 2$$

$$\leq = 4 \cdot 2$$

Date: \_\_\_\_\_

$$H(S) = 1$$

Information gain of B

$$\text{var } B >= 3$$

$$H(\text{var } B >= 3) = -\left[\frac{8}{12} \log_2 \frac{8}{12} + \frac{4}{12} \log_2 \frac{4}{12}\right]$$

$$H(\text{var } B >= 3) = 0.9182$$

$$\text{var } B < 3$$

$$H(\text{var } B < 3) = -\left[\frac{0}{4} \log_2 \frac{0}{4} - \frac{4}{4} \log_2 \frac{4}{4}\right]$$

$$H(\text{var } B < 3) = 0$$

$$IG_1(B) = 1 - \frac{12}{16}(0.9182) - \frac{4}{16}(0)$$

$$IG_1(B) = 0.7070$$

$$\text{Similarly } IG(C) = 0.54880$$

$$IG(D) = 0.41189$$

Information gain of C

$$\text{var } C >= 4 \cdot 2$$

$$H(\text{var } C >= 4 \cdot 2) = -\left(\frac{0}{6} \log_2 \frac{0}{6} + \frac{6}{6} \log_2 \frac{6}{6}\right)$$

$$H(\text{var } C >= 4 \cdot 2) = 0$$

$$H(\text{var } C <= 4 \cdot 2) = -\left(\frac{8}{10} \log_2 \frac{8}{10} + \frac{2}{10} \log_2 \frac{2}{10}\right)$$

$$H(\text{var } C <= 4 \cdot 2) = 0.72192$$

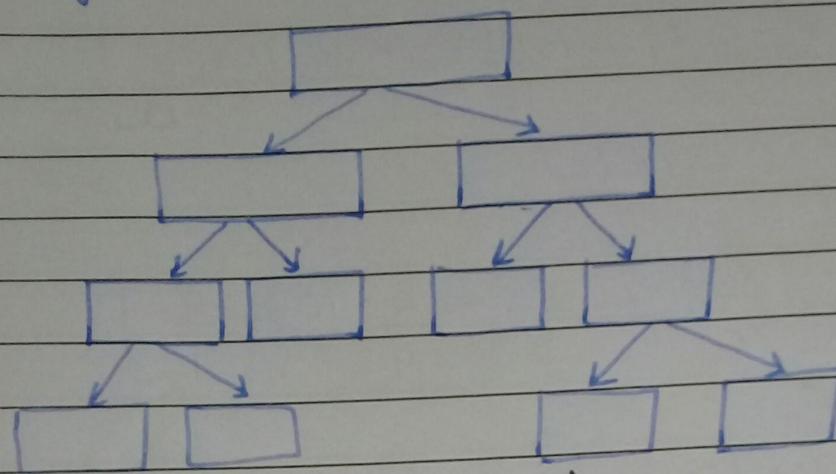
$$IG_1(C) = 1 - \left[\frac{6}{16} \times 0 + \frac{10}{16} \times 0.72192\right]$$

$$IG_1(C) = 0.4512$$

$$IG_1(C) = 0.5488$$

## RANDOM FOREST:

Building & Evaluating the random forest



1- Bagging (Random Forest)

2- Boosting (AdaBoost, Extreme Gradient Boosting)

RSD → Row Selection Data

Row Sample with replacement →

We can re-use data which we copied

→ We use mean as an alternate of Majority Voting  
in case of continuous Values

→ Random forest is a technique like decision tree

→ ensemble method

- جس کو چند trees میں سارے trees بنانے کا کہا جاتا ہے

→ Ensemble method uses two techniques

1- Bagging (Bootstrap aggregation) {Random Forest}

2- Boosting (Ada Boost, Extreme Gradient Boosting)

We are going to discuss bagging technique of Random Forest

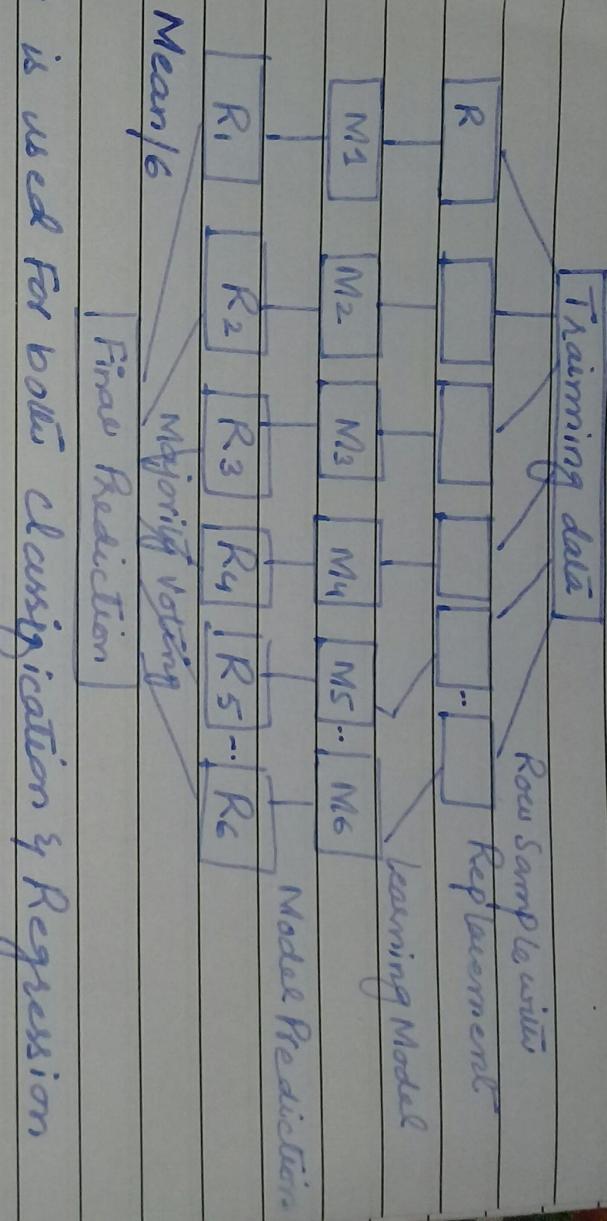
Ensemble → means → base

Test data sample

Bagging → ~~skew~~

Date: \_\_\_\_\_

Random Forest - Bagging



Bagging is used for both classification & Regression

$$\text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

$$SE = 139 / 139 + 32 = 0.8128 = 81\%$$

Cross Validation

Problem overcome by cross Validation

The biasness & ambiguity is removed through it

Holdout method ('a single test took out')

Kfold CV

Leave one out CV

Bootstrap Method

	Has Heart Disease	Does Not
Has Heart Disease	True Positive 139	False Positive 20
Doesn't have Heart D	False Negative 32	True Negative 112

$$\text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Rows → corresponds to what was predicted

Columns → corresponds to known truths

- \* True positives are patients that had heart disease that were also predicted to have heart disease.
  - \* True Negatives are patients that did not have heart disease & were predicted not to have heart disease.
  - \* False Negatives are when a patient has heart disease but the prediction says they did not
  - \* False Positive are patients that do not have heart disease but the prediction says that they do
- We can also calculate Sensitivity & Specificity through Confusion Matrix

$$SP = \frac{112}{112 + 20} = 0.8484 = 84.84\%$$

	Troll12	Gore Police	Cool as ice
Troll12	12	102	93 → FP
Gore Police	112	23	77 TN
Cool as ice	(83)	92	17

## Sensitivity (Troll12)

For Troll12: True Positive Troll12

True Positive Troll12 + False Negatives Troll12

$$S = \frac{12}{12 + 195} = 0.0579 = 5.79\%$$

## Sensitivity (Gore Police)

$$Se = \frac{23}{23 + 194}$$

$$= 0.11 = 11\%$$

## Sensitivity (Cool as ice)

$$Se = \frac{93}{93 + 94} = 0.4973 = 49.73\%$$

## Specificity (Troll12)

For Troll12:- True Negatives Troll12

True Negatives Troll12 + False Positives Troll12

$$SP = \frac{209}{209 + 195} = 0.5173 = 51.7\%$$

$$SP = \frac{205}{205 + 189} = 0.52 = 52\%$$

$$Se = \frac{17}{17 + 170} = 0.09$$

$$SP = \frac{249}{249 + 175} = 0.59$$

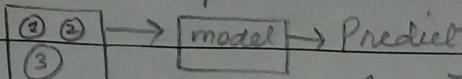
## Clustering

→ Unsupervised Learning Introduction:

If label then called  
Supervised learning

If donot label then  
First clustered & then

ML

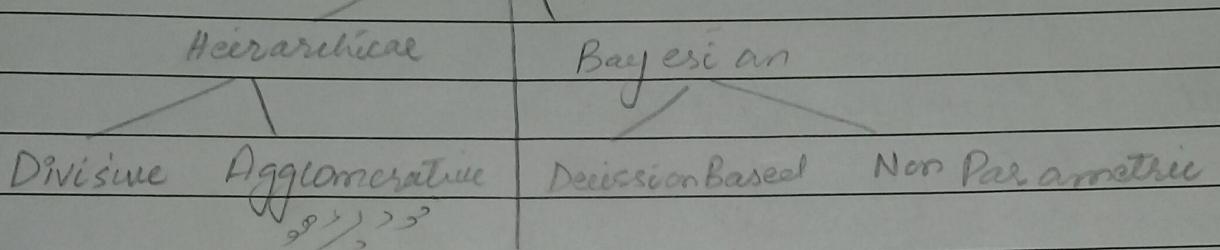


Cluster = same thing

unlabelled → means target is not known

clustering

1.5	1.5 - 4.5
4.6	4.6 - 2.5
7.1	7.1 - 15.9



Partitional

{Centroid} Model Based Graph theoretic Spectral

We are going to use centroid

Euclidean distance  $d(x_i, x_j) = \sqrt{\sum_{k=1}^d (x_i^{(k)} - x_j^{(k)})^2}$ 

clustering → is not only method

In centroid method find:

calculate  
X

→ Find Initial Seed (Randomly Picked)

→ Calculate distance from initial seed 1 &amp; seed 2

→ If the distance is less with S1 then belong to S1

When you got first cluster then you took out Mean / Mode &amp; then New Centroid

will be found then, the process repeat iteratively.

Calculate sum of squared error

Mean  
Mode  
X X Y Y

$\frac{10}{150} \times 100 =$	$\frac{800}{6.67}$	$A_1 A_2 \dots A_n$	TV
$x_1$	—	—	↓
$x_2$	—	0	Model $\Rightarrow$ <input type="checkbox"/>
$x_n$	—	:	Date: $\downarrow$ Y/No

correct  $140/150 = 93.3\%$

Variance matrix

$$CM = \mathbf{a}^T \mathbf{a} = \begin{bmatrix} 30 & 30 & 0 & -30 & -30 \\ 30 & 10 & 0 & -10 & -30 \\ -30 & 10 & 0 & 0 & 20 \end{bmatrix} \begin{bmatrix} 30 & 30 & -30 \\ 30 & 10 & 10 \\ 0 & 0 & 0 \\ -30 & -10 & 0 \\ -30 & -30 & 20 \end{bmatrix}$$

\* Same accuracy will be obtain if we reduce dataset.

$$CM = \mathbf{a}^T \mathbf{a} = \begin{bmatrix} 3600 & 2400 & -1200 \\ 2400 & 2000 & -1400 \\ -1300 & -1400 & 1400 \end{bmatrix} = \begin{bmatrix} MM & MP & ME \\ PM & PP & PE \\ EM & EP & EE \end{bmatrix}$$

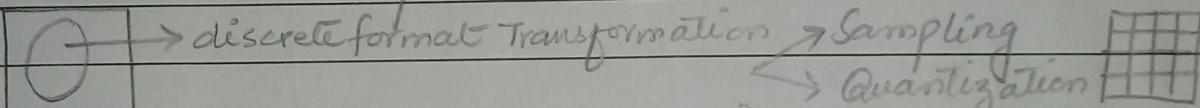
+ve: directly proportional

-ve: Inversely proportional

1- How to reduce the attributes?

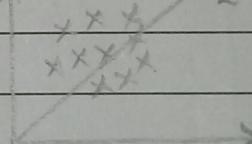
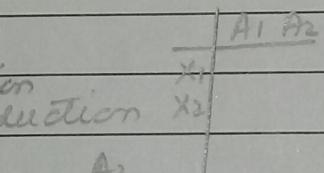
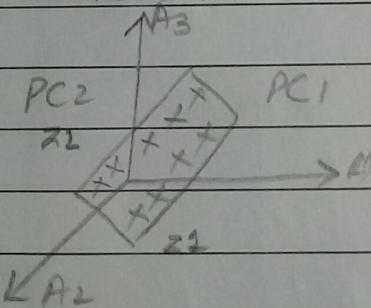
2- Whenever attributes will be reduced leads to effect the accuracy.

3- This also used as feature extraction



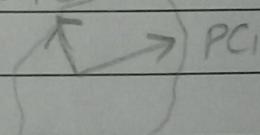
Feature vector

Feature extraction  
→ dimension reduction



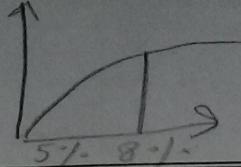
Eigen spaces  $PC_{35}$

Cartesian space



$PC_{35}$

9000 x 5



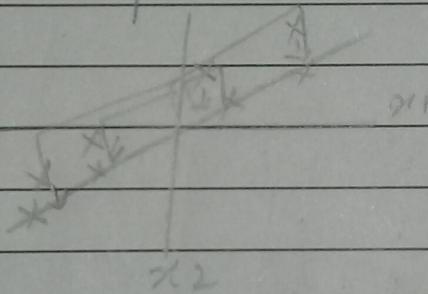
Date: \_\_\_\_\_

- How can we move from one space to another?
- The co-ordinates will be same after moving to other space?
- Yes, the co-ordinates will remain same.
- In those, 9000 which are important PC's?
- How to pick PC's?

\* Eigen space is a collection of principle components  
\* We used Principal Component Analysis for Entering into a new space.

PC<sub>1</sub>

Orthogonal  $\rightarrow$  independant  $\&$   $90^\circ$  angle

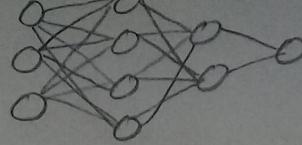


Plot the attribute on the line at the angle of  $90^\circ$ .  
Then take distance for each ' $x$ ' to the origin  
then plot different PC's for finding best PC<sub>1</sub>

Highest Distance = PC<sub>1</sub>

Technical Word = Map Honay k bad jiska  
Variance ziada Hoga, wo pick Hoga

# ML Lecture



Date: \_\_\_\_\_

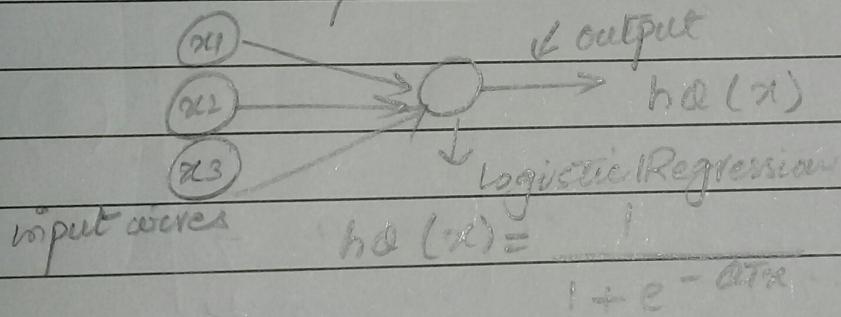
- \* Neural Networks learning it's own features
- Can we get a one unit neural network to compute this goal

Input  $Y \rightarrow$  Dendrite      |  
 Output  $Y \rightarrow$  Axon      | cell of Body  
 Nucleus

→ features receiving from a source  
 → we generate layers

- Neuron perform specific tasks
- Preprocessing is done through the layers
- Neuron is also called "Active Function"

- \* Feed input via input wires
- \* Logistic unit does computation
- \* Sends output down



Edges are called "Weight"

Date: \_\_\_\_\_

## Covariance of Multi Dimensional Model

- Covariance produces the relationship between two or more random variables.

Suppose

$$A = \begin{bmatrix} 90 & 80 & 40 \\ 90 & 60 & 80 \\ 60 & 50 & 70 \\ 30 & 40 & 70 \\ 30 & 20 & 90 \end{bmatrix}$$

Score of 3 subjects  
S - Students

$$\left[ \begin{array}{ccc} \frac{300}{300/5} & \frac{250}{250/5} & \frac{350}{350/5} \end{array} \right] \text{Total size}$$

$$\left[ \begin{array}{ccc} 60 & 50 & 70 \end{array} \right] \text{Mean score}$$

To calculate the deviation matrix  $a$

$$a = A - [1] A \cdot \left( \frac{1}{5} \right)$$

$\downarrow$  unit matrix       $\sim$  unit matrix

$$a = \begin{bmatrix} 90 & 80 & 40 \\ 90 & 60 & 80 \\ 60 & 50 & 70 \\ 30 & 40 & 70 \\ 30 & 20 & 90 \end{bmatrix} - \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix} \cdot \left( \frac{1}{5} \right)$$

$$a = \begin{bmatrix} 90 & 80 & 40 \\ 90 & 60 & 80 \\ 60 & 50 & 70 \\ 30 & 40 & 70 \\ 30 & 20 & 90 \end{bmatrix} - \begin{bmatrix} 50 & 50 & 70 \\ 60 & 50 & 70 \\ 60 & 50 & 70 \\ 60 & 50 & 70 \\ 80 & 50 & 70 \end{bmatrix}$$

$$a = \begin{bmatrix} 30 & 30 & -30 \\ 30 & 10 & 10 \\ 0 & 0 & 0 \\ -30 & -10 & 0 \\ -30 & -30 & 20 \end{bmatrix}$$

-ve: Inversely Proportional

$$\sqrt{(x_i - \bar{x})^2 + (y_i - \bar{y})^2}$$

Date: \_\_\_\_\_

+ve directly proportional

Co-Variance Matrix

$$CM = a^T \cdot a = \begin{bmatrix} 30 & 30 & 0 & -30 & -30 \\ 30 & 10 & 0 & -10 & -30 \\ -30 & 10 & 0 & 0 & 20 \end{bmatrix} \begin{bmatrix} 30 & 30 & -30 \\ 30 & 10 & 10 \\ 0 & 0 & 0 \\ -30 & -10 & 0 \\ -30 & -30 & 20 \end{bmatrix}$$

$$CM = \begin{bmatrix} 3600 & 2400 & -1200 \\ 2400 & 2000 & -1400 \\ -1200 & -1400 & 1400 \end{bmatrix} = \begin{bmatrix} MM & MP & ME \\ PM & PP & RE \\ EM & ER & EE \end{bmatrix}$$

$$\begin{bmatrix} 30 \times 30 + 30 \times 30 + 0 \times 0 + -30 \times -30 + -30 \times -30 \\ 30 \times 30 + 10 \times 10 + 0 \times 0 + -10 \times -10 + -30 \times -30 \end{bmatrix}$$