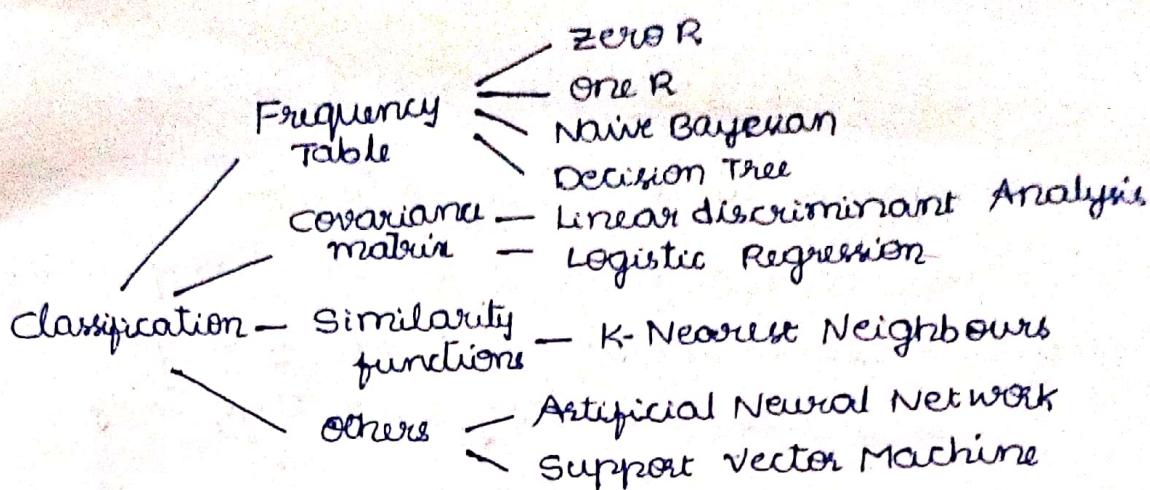


UNIT-4

Classification Methods :-



Classification:- It is the form of data analysis that extracts models describing important data classes.

Ex:- A class has A, B, C, D, E classes based on their Grades. By using a model, we can classify the students depends upon grades.

90 \geq grade class A
80 \geq , $<$ 90 class B
70 \geq , $<$ 80 C
60 \geq , $<$ 70 D
 $<$ 60 E

Applications of classification :-

- Image & pattern Matchings
- Medical diagnosis
- fraud detections
- Loan Approvals
- classifying financial market trends

Classification problem :-

Given a db $D = T_1, T_2 \dots T_n$ of Tuples (items, records) and a set of classes $C = C_1, C_2 \dots C_m$

The classification problem is to define mapping $f: D \rightarrow C$ where each T_i is assigned to one class. A class C_j contains precisely those tuples mapped to it.

$$G = \{t_i \mid f(t_i) = c_j; i \leq i \leq n \text{ & } t_i \in D\}$$

General Approach to the classification :-

There are two ways to do the classification process - Learning stage (or) Training stage
- classification stage (or) Testing Stage

Learning stage :- where a classification model is constructed

Classification stage :- where the model is used to predict the class labels for given data.

Ex:-

Training data set

Name	age	income	loan-decision
A	Youth	low	risk
B	Youth	middle	risk
C	middle	high	Safe
D	middle	low	risk
E	senior	low	risk
F	senior	low	risk

Stage 1
Learning stage

↓
classification Algo

↓
classification rules

↓

If age = youth then loan-decision = risk

If age = middle and income = high then loan decision = safe

If age = middle and income = low then loan decision = risk

If age = senior then loan-decision = risk



Issues in classification :-

- Missing data
- performance Measurement

Missing data :-

- . The missing values in the dataset create a problem during both training & classification problem itself.
- . The missing values in training data must be handled otherwise produced inaccurate results.
- . By using handling methods of missing data, we can overcome the problem.

Measuring the performance :-

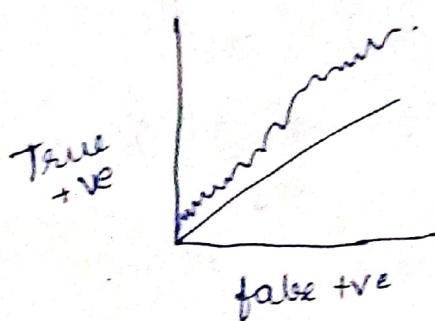
- The performance of a classification algo is usually examined by evaluating the accuracy of classifier.
- To determine which is best algo depends on the interpretation of problem by user.

- The space & time complexity is used to evaluate the algorithm but these approaches are usually secondary.

* By using the confusion matrix we can estimate the accuracy of classification methods.

x/ci	Yes	No
Class A	True +ve	false -ve
Class B	false +ve	True -ve

ROC curve:- The ROC curve (Receiver operating curve) shows the relationship b/w false +ve & true +ve.



Statistical Based Algorithm:-

1. Regression :- 1st unit refer
2. Bayesian method / Baye's classification

Naive Baye's classification method :-

It is a statistical classifier they can predict the class membership probabilities such as the probability that a given tuple belongs to the particular class.

The Bayesian classification is based on baye's

Theorem.

→ Naive Bayesian classifier assume that the Effect of an attribute value on a given class is independent of the value of other attributes. This assumption is called as class conditional Independence.

Baye's Theorem:-

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

where $P(A|B)$ = The probability of occurrence of event A given the event B is true.

$P(A) \& P(B)$ = The probability of occurrence of Event A and Event B

A is called proposition

B is called evidence

$P(A)$ is called prior probability of proposition

$P(B)$ is called prior probability of evidence

$P(A|B)$ is called posterior or posterior probability.

$P(B|A)$ is called likelihood.

$$\text{posterior} = \frac{(\text{likelihood})(\text{prior probability of proposition})}{\text{prior probability of evidence}}$$

Ex:- consider the given dataset apply the Naive bayesian algo and predict that if a fruit has the following properties then which type of fruit it is.

fruit = {yellow, sweet, long}

Fruit	Yellow	Sweet	long	Total
Mango	350	450	0	650
Banana	400	300	350	400
Other	50	100	50	150
Total	800	850	400	1200

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

$P(X/\text{mango})$

$P(\text{Yellow}/\text{mango})$

$$= \frac{P(m/y) P(y)}{P(m)}$$

$$= \frac{350}{800} \times \frac{800}{1200}$$

$$\frac{650}{1200}$$

$$= 0.52.$$

$P(\text{Sweet}/\text{mango})$

$$= \frac{P(\text{mango/sweet}) P(\text{sweet})}{P(\text{mango})}$$

$$= \frac{450}{850} \times \frac{800}{1200}$$

$$\frac{650}{1200}$$

$$= 0.69$$

$P(\text{long}/\text{mango})$

$$= \frac{P(\text{mango/long}) P(\text{long})}{P(\text{mango})}$$

$$= \frac{0}{400} \times \frac{400}{1200} = 0$$

$$\frac{650}{1200}$$

Result: $\pi p(x_i | G_i)$

$$= 0.52 \times 0.69 \times 0$$

$$= 0$$

Ans: Banana

$P(X/\text{Banana})$

$$P(Y|B) = \frac{P(B|Y) P(Y)}{P(B)}$$

$$= \frac{\frac{400}{800} \times \frac{800}{1200}}{\frac{400}{1200}} = 1$$

$P(\text{Sweet}/\text{Banana})$

$$\frac{P(B|S) P(S)}{P(B)}$$

$$= \frac{300}{850} \times \frac{850}{1200}$$

$$\frac{400}{1200}$$

$$= 0.75$$

$P(\text{long}/\text{Banana})$

$$\frac{P(B|L) P(L)}{P(B)}$$

$$= \frac{350}{400} \times \frac{400}{1200}$$

$$\frac{400}{1200}$$

$$= 0.87$$

Result: 0.75×0.87

$$= 0.65$$

As it is highest
Banana

$P(X/\text{others})$

$$P(Y|O) = \frac{P(O|Y) P(Y)}{P(O)}$$

$$= \frac{\frac{50}{800} \times \frac{800}{1200}}{\frac{150}{1200}} = 0.33$$

$$P(S|O) = 0.66$$

$$P(L|O) = 0.33$$

$$\text{result} = 0.072$$

Ex-2

ID.	age	income	Student	Credit-rating	Class: buys-computer
1	youth	high	NO	fair	NO
2	y	high	NO	Excellent	NO
3	middle	high	NO	f	yes
4	senior	medium	y	f	y
5	s	low	y	f	NO
6	s	low	y	ex	y
7	m	low	N	ex	NO
8	y	medium	y	f	y
9	y	low	y	f	NO
10	s	m	y	f	y
11	y	m	N	ex	y
12	m	m	y	ex	y
13	m	h	N	f	NO
14	s	m	N	ex	

$x = (\underline{\text{age} = \text{youth}}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$

$$P(\text{buys-computer} = \text{Yes}) = 9/14$$

$$P(\text{buys-computer} = \text{NO}) = 5/14$$

$$P(x / \text{buys-computer} = \text{Yes})$$

$$P(\text{age} = \text{youth} / \text{buys-computer} = \text{Yes}) = 2/9$$

$$P(\text{age} = \text{youth} / \text{buys-computer} = \text{NO}) = 3/5$$

$$P(\text{income} = \text{medium} / \text{buys-computer} = \text{Y}) = 4/9$$

$$P(\text{income} = \text{medium} / \text{buys-computer} = \text{N}) = 2/5$$

$$P(\text{income} = \text{medium} / \text{buys-computer} = \text{Y}) = 6/9$$

$$P(\text{student} = \text{yes} / \text{buys comp} = \text{NO}) = 2/5$$

$$P(\text{st} = \text{Y} / \text{buys comp} = \text{NO}) = \frac{2}{5} \text{ total no. of NO's}$$

$$P(g_1 = \text{fair} / \text{bc} = \text{Y}) = 6/9$$

$$P(g_2 = \text{fair} / \text{bc} = \text{N}) = 2/5$$

$$\text{Yes} \rightarrow \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9} \times \frac{9}{14} = 0.044 \quad \therefore x \text{ will buy}$$

$$\text{NO} \rightarrow \frac{3}{14} \times \frac{2}{5} \times \frac{2}{5} \times \frac{3}{5} \times \frac{5}{14} = 0.038.$$

the computer

Decision tree Induction:-

A decision tree is a flowchart like tree structure where each internal node (non leaf node) denotes a test on an attribute. Each branch represents an outcome of the test and the leaf node represents a class label. The top most node in tree is the root node.

In this decision tree a tuple x for which class label is unknown, attribute values of the tuple are tested against the decision tree. A path is traced from root to the leaf node which holds the class prediction for the tuple.

In this decision tree does not require any domain knowledge or parameter setting to classify the data (so it is very popular)

Applications :-

- Manufacturing & production
- financial Analysis
- Astronomy
- Molecular Biology.

KNN → K-Nearast Neighbours

The KNN classifier to predict the target label by finding nearest neighbour class.

The closest class will be identified using the distance measures like Euclidean distance or Manhattan distance or Taxicab distance.

Application:-

- used in Recommendation systems
- To store info in servers.

How to choose the K-value:-

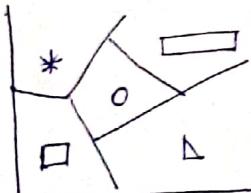
(Select the) Selecting the K-value in KNN is the most critical problem.

- . The simple approach to select the K-value is \sqrt{n} where n is no. of observations / data points
- . By using the cross validation also, we can select the value of K.

NOTE :-

The KNN Algorithm also called as lazy learner algo because we are not used to learn any function or to store values.

Case1 :-



Called as voronoefft
as all classes are nearer to the obj we are unable to select the class.

Perform KNN Algo on following dataset & predict the class for $x(P_1=3 \text{ & } P_2=7)$ & $K=3$

	P_1	P_2	classes
①	7	7	False
②	7	4	False
③	3	4	True
④	4	4	True

Euclidean distance

$$\sqrt{(x_H - H)^2 + (y_H - H)^2}$$

/ \

 Observed Actual
 or given value value

Given $P_1 = 3, P_2 = 7$

As $K = 3$

consider 3 values
from low to high

$$① \sqrt{(3-7)^2 + (7-7)^2} = \sqrt{16} = 4 - \text{True } \checkmark$$

$$② \sqrt{(3-7)^2 + (7-4)^2} = \sqrt{25} = 5$$

$$③ \sqrt{(3-3)^2 + (7-4)^2} = \sqrt{13} = 3.6 - \text{False}$$

$$④ \sqrt{(3-1)^2 + (7-4)^2} = \sqrt{13} = 3.6 - \text{True } \checkmark \quad \text{highest category}$$

Case 2 :-

Line encoding

Gender

female - 0	}	consider one as zero's & other to one's
male - 1		
male - 1		
female - 0		

Range Encoding

	Height (cm)	Avg	Lower	Upper
consider Avg	100-110	105	100	110
	110-120	115	110	120
	120-130	125	120	130
	130-140	135	130	140

Case 3 :-

	Age	Age youth	Age middle	Age senior
divide them	youth	1	0	0
	middle age	0	1	0
	senior	0	0	1
	youth	1	0	0
	middle age	0	1	0
	senior	0	0	1
		0	0	1

Attribute Selection :- (Decision Tree)

An attribute selection is a heuristic for selecting the decision tree splitting criteria that best separate the given data partition (D) of class label training tuples into individual classes.

It is also called as splitting rules because they determine how the tuples of given node are to be split.

- The attribute selection provides a ranking for each attribute and select the best score attribute as a root node of trees.
- There are three popular attribute selection methods.
 - Information gain
 - Gain ratio
 - Gini Index

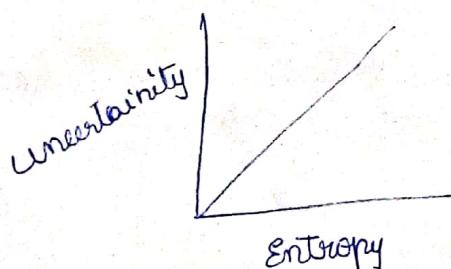
$$\text{Information gain } (D, A) = \text{entropy}(D) - \sum_{j=1}^{|D_j|} \frac{|D_j|}{|D|} \text{entropy}(D_j)$$

$$\text{Entropy } (D) = \sum_{i=1}^C -p_i \log_2 (p_i)$$

$$= -\frac{p_i}{p+n} \log_2 \left(\frac{p_i}{p+n} \right) - \frac{n_i}{p+n} \log_2 \left(\frac{n_i}{n+p} \right)$$

- * Information gain is used as Attribute selection measure
- * Entropy is used to find out uncertainty associated with a random variable.

The values of an entropy is in b/w 0 to 1.



From previous example

RID	age	income	student	credit rating	buys computer
-----	-----	--------	---------	---------------	---------------

$$P_i \rightarrow \text{Yes} \Rightarrow 9$$

$$n_i \rightarrow \text{No} \Rightarrow 5$$

$$\begin{aligned}
 \text{entropy (D)} &= \frac{-P_i}{P+n} \log_2 \left(\frac{P_i}{P+n} \right) - \frac{n_i}{P+n} \log_2 \left(\frac{n_i}{P+n} \right) \\
 &= -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) \\
 &= -\frac{9}{14} \log_2 (0.6428) - \frac{5}{14} \log_2 (0.3571) \\
 &= -0.6428 \times \frac{\log(0.6428)}{\log 2} - \frac{5}{14} \times \frac{\log(0.3571)}{\log 2} \\
 &= -0.6428 \times \frac{-0.1919}{0.3010} - 0.3571 \times \frac{0.4472}{0.3010} \\
 &= -0.6428 \times -0.6375 - 0.3571 \times -1.4857 \\
 &= 0.4097 + 0.5305 \\
 &= 0.9402
 \end{aligned}$$

Info age (D)	youth	middle	senior
	Yes = 2	Yes = 4	Yes = 3
	No = 3	No = 0	No = 2

$$\begin{aligned}
 &= \frac{5}{14} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{4}{14} \left(-\frac{4}{4} \log_2 \frac{4}{4} - 0 \right) \\
 &\quad + \frac{5}{14} \left(\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\
 &= 0.694 \text{ bits}
 \end{aligned}$$

$$\text{Information gain (D, age)} = \text{info}(D) - \text{info.age}(D)$$

$$= 0.940 - 0.694$$

$$= 0.246 \text{ bits}$$

Info Income (D) High medium low

Yes = 2	Yes = 4	Yes = 3
No = 2	No = 2	No = 1

$$\begin{aligned}
 &= \frac{4}{14} \left(-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right) + \frac{6}{14} \left(-\frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6} \right) \\
 &\quad + \frac{4}{14} \left(\frac{2}{4} \log \frac{2}{4} - \frac{1}{4} \log \frac{1}{4} \right) \\
 &= 0.909
 \end{aligned}$$

Information gain(D, Income) = Info(D) - InfoIncome(D)

$$= 0.940 - 0.909$$

$$= 0.029 \text{ bits}$$

Info Student (D) Yes No

Yes = 6	Yes = 3
No = 1	No = 4

$$= 0.786$$

Info gain = 0.151 bits.

Credit rating fair excellent

Y = 6	Yes = 3
N = 2	No = 3

~~$= 0.8926$~~

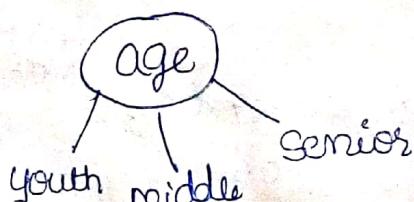
~~$= 0.940 - 0.8926 = 0.048$~~

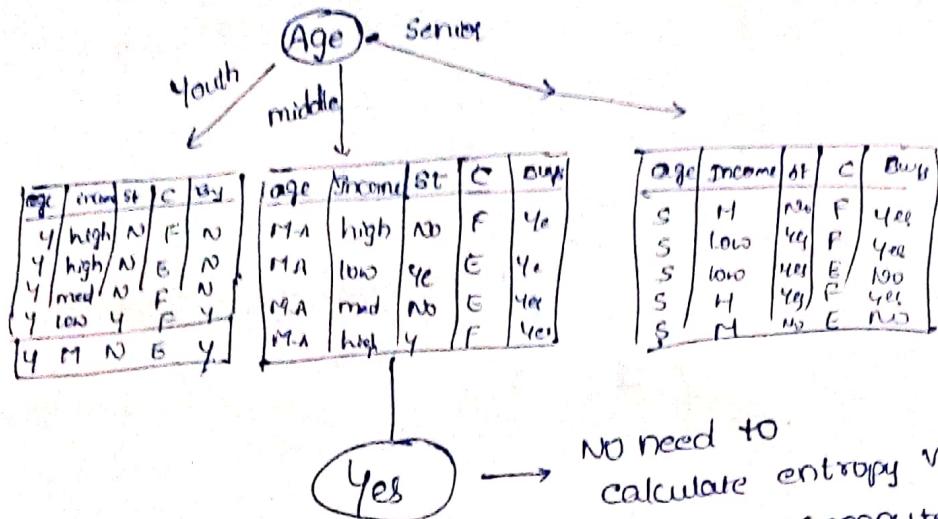
Gain Y, Age = 0.246 ✓

Gain Insurance = 0.029

Student = 0.154

Credit = 0.048





Yes

No need to calculate entropy values
bcz Buys-computer = Yes
on all tuples.

$$\text{Entropy Youth (D)} = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5}$$

$$\begin{aligned} \text{NO} &= 3 \\ \text{Yes} &= 2 \end{aligned} \quad = 0.2922$$

Next attribute	High	Medium	Low
Income	Yes = 0 NO = 2	Yes = 2 NO = 1	Yes = 1 NO = 0

$$\begin{aligned} &= \frac{2}{5} \left(-\frac{0}{2} \log \frac{0}{2} - \frac{2}{2} \log \frac{2}{2} \right) + \frac{2}{5} \left(-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right) + \frac{1}{5} \left(-\frac{1}{1} \log \frac{1}{1} - 0 \right) \\ &= \frac{2}{5} (-\log \frac{1}{2}) = 0.1204 \end{aligned}$$

$$\text{Information gain} = 0.29 - 0.1204 = 0.17$$

Entropy Student (D)	Yes	NO
	Yes = 2 NO = 0	Yes = 0 NO = 3

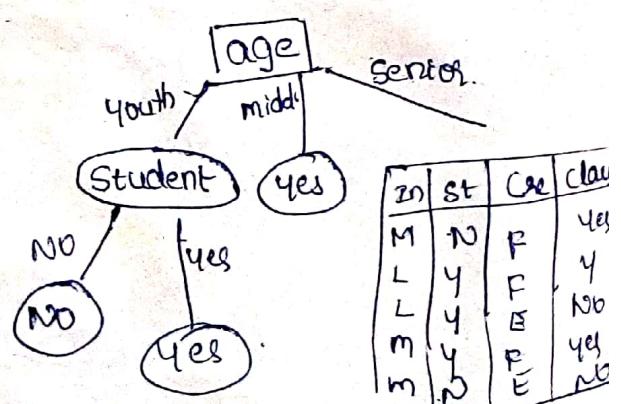
$$\frac{2}{5} \left(-\frac{2}{2} \log \left(\frac{2}{2} \right) - \frac{3}{2} \log \left(\frac{3}{2} \right) \right) + \frac{3}{5} \left(-\frac{0}{3} \log \frac{0}{3} - \frac{3}{3} \log \frac{3}{3} \right) = 0$$

$$\text{Info gain} = 0.29 - 0 = 0.29$$

Entropy Credit rating (D)	Fair	Excellent	
	Yes = 1 NO = 2	Yes = 1 NO = 1	= 0.95

$$\begin{aligned} \text{Info gain} &= 0.29 - 0.95 \\ &= -0.66 \end{aligned}$$

Among all student have high information gain so next node is student



Entropy senior (D) $N_{\text{No}}=2$ $N_{\text{Yes}}=3$

$$-\frac{3}{5} \log\left(\frac{3}{5}\right) - 2/5 \log\left(\frac{2}{5}\right) = 0.29$$

Entropy Income (D) Medium Low

$N_{\text{Yes}}=2$	$N_{\text{No}}=1$
$N_{\text{No}}=1$	$N_{\text{Yes}}=1$

$$\frac{3}{5} \left(-\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) \right) + \frac{2}{5} \left(-\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) \right) = 0.286$$

$$\text{Info gain (Income)} = 0.29 - 0.286 = 0.006.$$

Entropy credit rating (D) Fair Excellent.

$N_{\text{Yes}}=3$	$N_{\text{No}}=0$
$N_{\text{No}}=0$	$N_{\text{Yes}}=2$

$$\frac{3}{5} \left(-\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) \right) + \frac{2}{5} \left(0/2 \log 0/2 - \frac{1}{2} \log\left(\frac{1}{2}\right) \right) = 0$$

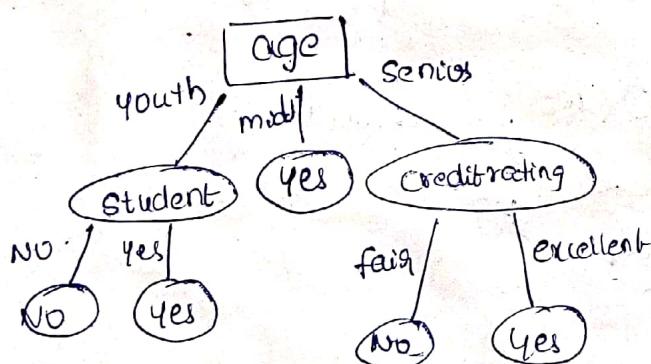
$$\text{Info gain (credit rating)} = 0.29 - 0 = 0.29$$

Entropy student (D) Yes No

$N_{\text{Yes}}=2$	$N_{\text{No}}=1$
$N_{\text{No}}=1$	$N_{\text{Yes}}=1$

$$\frac{3}{5} \left(-\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) \right) + \frac{2}{5} \left(-\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) \right) = 0.286$$

$$\text{Info gain} = 0.292 - 0.286 = 0.006$$



✓ Gain Ratio (C4.5), Gini Index

Refer

TB

Artificial Neural Networks :-

Neural networks are parallel computing device which is basically an attempt to make a computer model of the brain.

The main objective of neural networks to perform the various computation tasks faster than tradition system

The tasks are - pattern Recognition

- classification

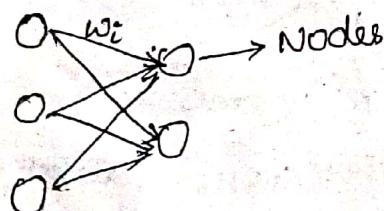
- Approximation

- clustering & optimization

* ANN also called as Artificial neural system, parallel distributed processing system. connectionist systems.

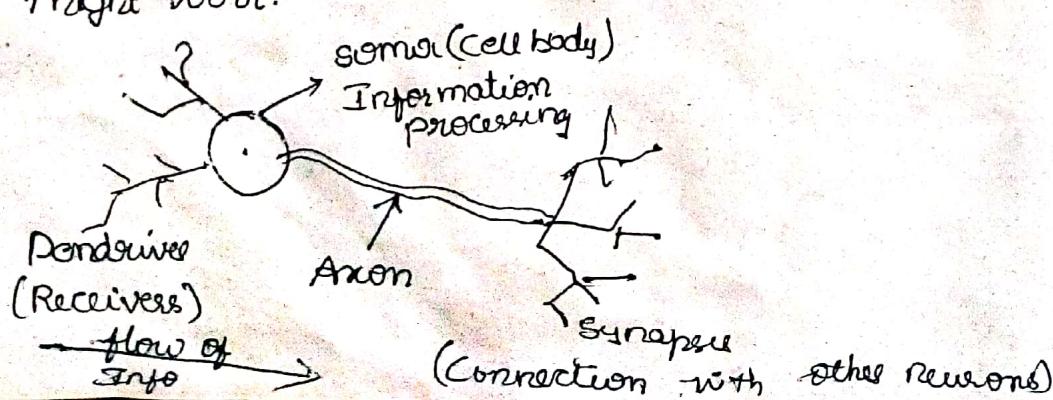
Working Of ANN:-

- * ANN requires the large collection of units that are interconnected in some pattern to allow communication b/w the units. Each unit is called as neuron/node. which are simple processors operates in parallel.
- * Each connection link is associated with weight that has information about the input signal



History of ANN:-

- In 1943, started the origin of networks using electric circuits in order to describe how neurons in brain might work.



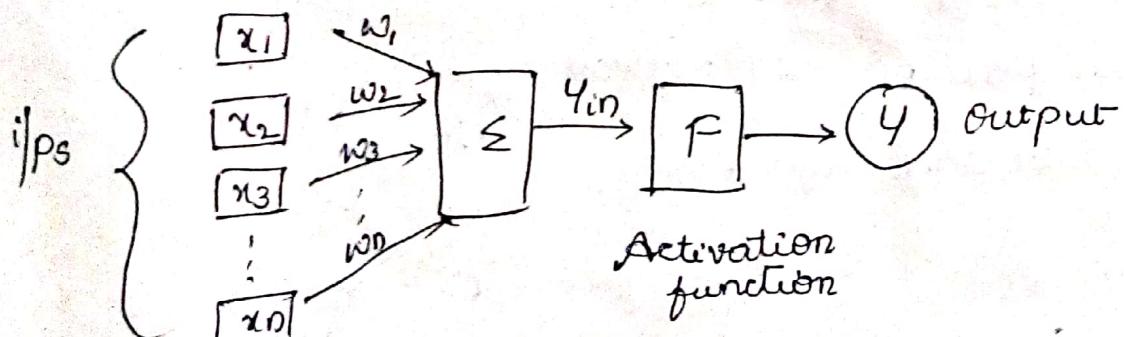
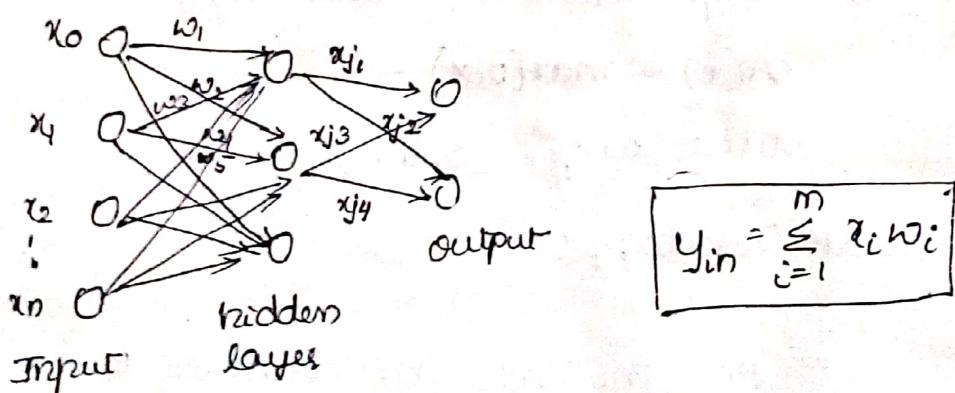
BNN
Soma
Dendrites
Synapse
Axon

ANN
Node
Input
weights / interconnection
output

The difference b/w BNN & ANN :-

Criteria	BNN	ANN
processing	massively parallel, slow but superior than ANN	massively parallel fast but inferior than BNN
Size	10^8 neurons & 10^{15} connections	10^2 to 10^4 nodes depends on application & new design
Fault Tolerance	performance degrades with even partial damage	performance is very fast

General Model of ANN :-



Activation Functions :-

- To check node is working or not
- convert the data into non-linear form

Types:-

- Step function
- Sigmoid function
- Linear function
- Tanh
- ReLU

Step Function : $A = 1$ if $y > \text{Threshold}$ otherwise $A = 0$

Sigmoid function :-

$$\frac{1}{1+e^{-x}}$$

Q types - binary 0 to 1
bipolar -1 to +1

Linear function : $y = mx$ ($-\infty$ to $+\infty$)

Tanh : Tangent hyperbolic function

$$\text{Tanh} = \frac{2}{1+e^{2x}} - 1 = Q * \text{sigmoid}(2x) - 1$$

ReLU (Rectified Linear unit)

It is implemented in Hidden layers

$$A(x) = \max(0, x)$$

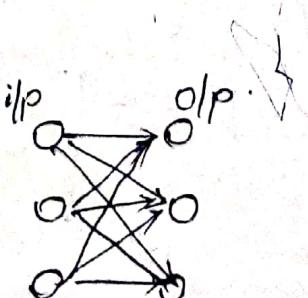
value range = [0, inf)

Network Topology :-

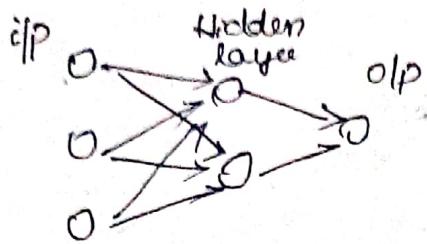
How the nodes are connected to each other.

- Feed forward Network

(1) Single layer F.F.N

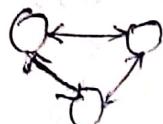


(ii) Multilayer FFN

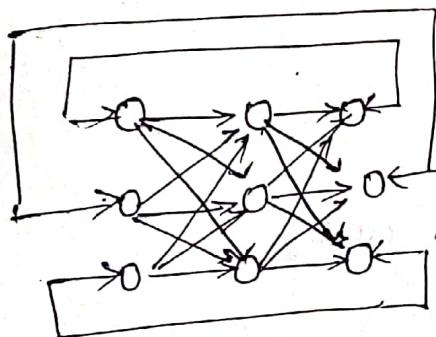


- Feed back N/w

(i) Fully recurrent N/w



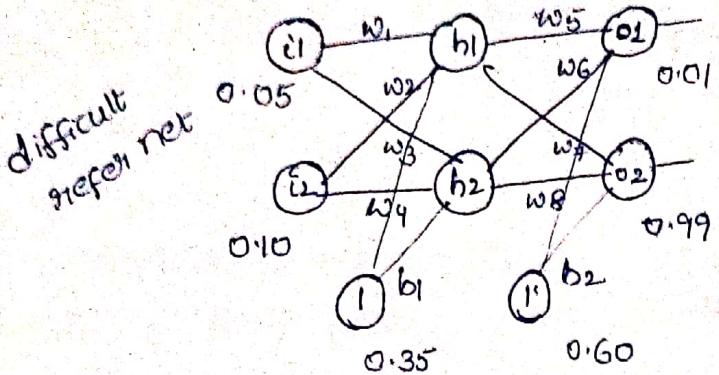
(ii) Jordan N/w



Issues in ANN :-

- (i) Attributes (No. of source nodes)
- (ii) No. of hidden layers
- (iii) No. of hidden nodes
- (iv) Training data
- (v) No. of sinks
- (vi) Interconnections
- (vii) weights
- (viii) Activation function
- (ix) Learning Techniques
- (x) Stop

Ex: Back propagation



$$\begin{aligned}
 w_1 &= 0.15 \\
 w_2 &= 0.20 \\
 w_3 &= 0.25 \\
 w_4 &= 0.30 \\
 w_5 &= 0.40 \\
 w_6 &= 0.45 \\
 w_7 &= 0.50 \\
 w_8 &= 0.55
 \end{aligned}$$

Back propagation :-

Algorithm -

- for each input layer unit j

$$o_j = I_j$$

- for each hidden or output layer unit j

$$I_j = \sum_i w_{ij} + \theta_j$$

Back propagation errors :-

for each unit j in the O/p layer

$$Err_j = o_j(1-o_j)(T_j - o_j)$$

for each unit j in the hidden layer

$$Err_j = o_j(1-o_j) \sum_k Err_k w_{jk}$$

for each weight w_{ij} in the N/W

$$\Delta w_{ij} = (\eta) Err_j o_j$$

$$w_{ij} = w_{ij} + \Delta w_{ij}$$

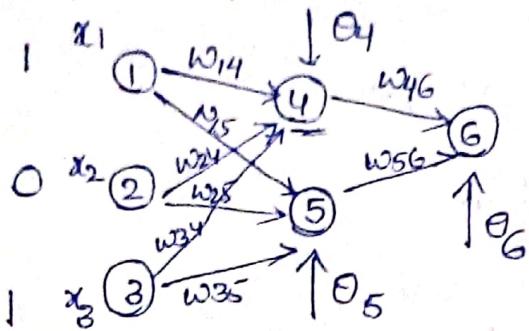
for each bias θ_j in N/W

$$\Delta \theta_j = (\eta) Err_j$$

$$\theta_j = \theta_j + \Delta \theta_j$$

[for each hidden or output layer unit j]

Ex-



we are needed
to get O_1 as 1
if not modify
the weights & modi

Given

$$x_1 \quad x_2 \quad x_3$$

$$1 \quad 0 \quad 1$$

$$\begin{aligned} w_{14} &= 0.2 & w_{34} &= 0.5 \\ w_{15} &= 0.3 & w_{35} &= 0.2 \\ w_{24} &= 0.4 & w_{46} &= -0.3 \\ w_{25} &= 0.1 & w_{56} &= -0.2 \end{aligned}$$

$$\theta_4 = -0.4$$

$$\theta_5 = 0.2$$

$$\theta_6 = 0.1$$

$$O_j = T_j$$

$$O_1 = x_1$$

$$O_2 = x_2$$

$$O_3 = x_3$$

unit j	input layer	output
4	$x_1 \times w_{14} + x_2 \times w_{24} + x_3 \times w_{34} + \theta_4$	$\frac{1}{1+e^{-x}} = \frac{1}{1+e^{-0.7}} = 0.332$
	$= 1 \times 0.2 + 0 \times 0.3 + (0.5)1 + (-0.4)$	
	$= -0.7$	

5	$x_1 \times w_{15} + x_2 \times w_{25} + x_3 \times w_{35} + \theta_5$	$\frac{1}{1+e^{-x}} = \frac{1}{1+e^{-0.7}} = 0.668$
	$= 1 \times 0.3 + 0 \times 0.1 + 1 \times 0.2 + 0.2$	
	$= 0.7$	

6	$O\text{utput}(4) \times w_{46} + O(5) \times w_{56} + \theta_6$	$\frac{1}{1+e^{-x}} = \frac{1}{1+e^{-0.13}} = 0.467$
	$= (-0.33)(-0.3) + (0.668)(-0.2) + 0.1$	
	$= 0.332(-0.3) + 0.668(-0.2) + 0.1$	
	$= -0.13 \neq 1$	

so, we need to calculate errors from back

$$\text{Error}_j = O_j(1-O_j)(T_j - O_j)$$

$$\text{Error}_6 = O_6(1-O_6)(T_6 - O_6) = 0.467(1-0.467)(1-0.467) = 0.1326$$

error for hidden layer \rightarrow one layer means eliminated as it has only one node below

$$err_j = o_j(1-o_j) \sum_k e_{kj} w_{jk}$$

$$err_4 = 0.332(1-0.332)$$

$$err_5 = o_5(1-o_5) err_6 w_{56}$$

$$= 0.668(1-0.668)(0.132)(-0.2)$$

$$= -0.0058$$

$$err_4 = o_4(1-o_4) err_6 w_{46}$$

$$= 0.332(1-0.332)(0.132)(-0.3)$$

$$= -0.0087$$

calculation for weights & bias (modification)

$\Delta w_{ij} = w_{ij} + (l) \overbrace{err_j o_i}^{\text{represents learning rate given in question}}$

$$w_{46} = w_{46} + 0.9 err_6 o_4$$

$$= -0.3 + (0.9)(0.132)(0.332)$$

$$= -0.261$$

$$w_{56} = w_{56} + 0.9 err_6 o_5 = -0.12$$

$$w_{14} = w_{14} + 0.9 err_4 o_1 = 0.192$$

$$w_{15} = w_{15} + 0.9 err_5 o_1 = 0.294$$

$$w_{24} = w_{24} + 0.9 err_4 o_2 = 0.4$$

$$w_{25} = w_{25} + 0.9 err_5 o_2 = 0.1$$

$$w_{34} = w_{34} + 0.9 err_4 o_3 = -0.507$$

$$w_{35} = w_{35} + 0.9 err_5 o_3 = 0.194$$

changing bias values

$$\theta_j = \theta_j + (l) err_j$$

$$\theta_4 = \theta_4 + 0.9 err_4$$

$$= -0.407$$

$$= -0.205$$

$$\theta_5 = \theta_5 + 0.9 err_5 = 0.194$$

$$\theta_6 = \theta_6 + 0.9 err_6$$

$$= 0.219$$

$$\begin{aligned}
 \text{Op}_4 &\Rightarrow w_{14}x_1 + w_{24}x_2 + w_{34}x_3 + \theta_4 \\
 &= 0.192(1) + 0 + (-0.507) + -0.407 \\
 &= -0.722 \quad \text{Op}_4 = \frac{1}{1+e^{(-0.722)}} = 0.262 \\
 \\
 5 &\Rightarrow w_{15}x_1 + w_{25}x_2 + w_{35}x_3 + \theta_5 \quad 0.56 \\
 &= 0.4(w_{15}) + 0.5(w_{25}) + 0.5(w_{35}) \\
 6 &\Rightarrow w_{16}x_1 + w_{26}x_2 + \theta_6 \quad \frac{1}{1+e^{(0.839)}} \\
 &= 0.32 + 0.219 = 0.539 \\
 &(0.32)(-0.261) + (0.56)(-0.12) + 0.219 \quad \text{Op}_5 : 0.5 \\
 &= 0.06
 \end{aligned}$$

Rule Based Classification :-

By using If-then, we make the rules to classify the data.

Ex:- If age = middle then buys computer = yes
 It will be work like Divide & conquer method
 There are two methods in rule-based classification
 - zero R
 - one R

The General syntax for the Rule is IF condition THEN Conclusion

<u>Antecedent</u> (or) precondition	<u>Rule consequent</u> (or) (condition) \Rightarrow y
--	---

To find our the rule is satisfied or not, we use coverage & accuracy to validate the rule.

$$\text{Coverage}(R) = \frac{n_{\text{covers}}}{n_{\text{tuples}}} \quad \text{no. of Tuples in data}$$

$$\text{Accuracy}(R) = \frac{n_{\text{correct}}}{n_{\text{covers}}}$$

Rules are constructed two ways

- By using Decision Tree
- By using Neural Networks

5. Clustering

- collection of similar objects contains strong relationship and more dissimilar compare to other cluster.
- Intra relation is stronger than inter relation in cluster.

K-Means clustering Algo:-

→ No. of clusters, we are going to dividing given large data sets.

Ex:- Divide the given sample data into two clusters using K-Means Algorithm.

Height (H)	Weight (W)		
185	72		consider any one as centroid
170	50		
168	60		
179	68		
182	72	C_1	Height weight
188	77		185 72
180	71	C_2	170 50
180	70		Centroid

$$\text{Euclidean distance} = \sqrt{(x_H - H_i)^2 + (x_W - W_i)^2}$$

/ \ Observed value / \ Centroid value

$$C_1 = \{1$$

$$C_2 = \{2, 3 \} \rightarrow \text{as } C_2 \text{ gets 2 values, centroid changes}$$

$$C_2 \left(\frac{170+168}{2}, \frac{60+50}{2} \right)$$

$$new(R_3) = \sqrt{(168-185)^2 + (60-72)^2} \\ = 20.8$$

/ \ C_1 C_2

$$C_2 \rightarrow \sqrt{(168-190)^2 + (60-50)^2} = 10.48$$

Least go to C_2 under 3 write

$$new(R_4): C_1 \rightarrow \sqrt{(179-185)^2 + (68-72)^2} = \sqrt{52} = 7.2$$

Least go to C_1

$$C_2 \rightarrow \sqrt{(179-169)^2 + (68-55)^2} = \sqrt{269} = 16.4$$

$C_1 \{1, 4\}$ $C_2 \{2, 3\}$

Centroid

$$C_1 \left(\frac{185+179}{2}, \frac{72+68}{2} \right) = (182, 70)$$

Row (R₅) $C_1 \rightarrow \sqrt{(182-182)^2 + (72-70)^2} = 2 \text{ least}$
 $C_2 \rightarrow \sqrt{(182-169)^2 + (72-55)^2}$

 $C_1 \{1, 4, 5\}$ $C_2 \{2, 3\}$

Centroid

$$\left(\frac{182+182}{2}, \frac{70+72}{2} \right) = (82, 71)$$

Row (R₆) $C_1 \rightarrow \sqrt{(188-182)^2 + (77-71)^2} = \sqrt{36+36} \text{ least}$
 $C_2 \rightarrow \sqrt{(188-169)^2 + (77-55)^2}$

Row (R₇) $C_1 \{1, 4, 5, 6\}$ $C_2 \{2, 3\}$

Centroid

$$\left(\frac{182+188}{2}, \frac{71+77}{2} \right) = (185, 74)$$

$$C_1 \rightarrow \sqrt{(180-185)^2 + (71-74)^2} = 5.83 \text{ least}$$

$$C_2 \rightarrow \sqrt{(180-169)^2 + (71-55)^2} = 19.41$$

 $C_1 \{1, 4, 5, 6, 7\}$ $C_2 \{2, 3\}$

Centroid

$$\left(\frac{185+180}{2}, \frac{74+71}{2} \right) = (182.5, 72.5)$$

Row (R₈)

$$C_1 \rightarrow \sqrt{(180-82.5)^2 + (70-72.5)^2} = 35 \text{ least}$$

$$C_2 \rightarrow \sqrt{(80-169)^2 + (70-55)^2} = 18.6$$

 $C_1 \{1, 4, 5, 6, 7, 8\}$ $C_2 \{2, 3\}$

Clustering : 1) partitional based 2) hierarchical

3) density based 4) clustering LB Agglomerative ✓ Divisive

- In hierarchical clustering method works by grouping data objects into hierarchy or tree of clusters
- There are different sub methods in hierarchical clustering.

1. Agglomerative Nesting vs divisive analysis
hierarchical clustering.

2. Distance measures in algorithmic methods

3. BIRCH (balanced iterative reducing & clustering using hierarchies) multiphase hierarchical clustering

4. chameleon multiphase hierarchical clustering

Agglomerative Nesting vs divisive :-

A hierarchical clustering method can be either agglomerative or divisive, depending on whether the hierarchical decomposition is formed in a bottom up (merging) or top down (splitting) approach.

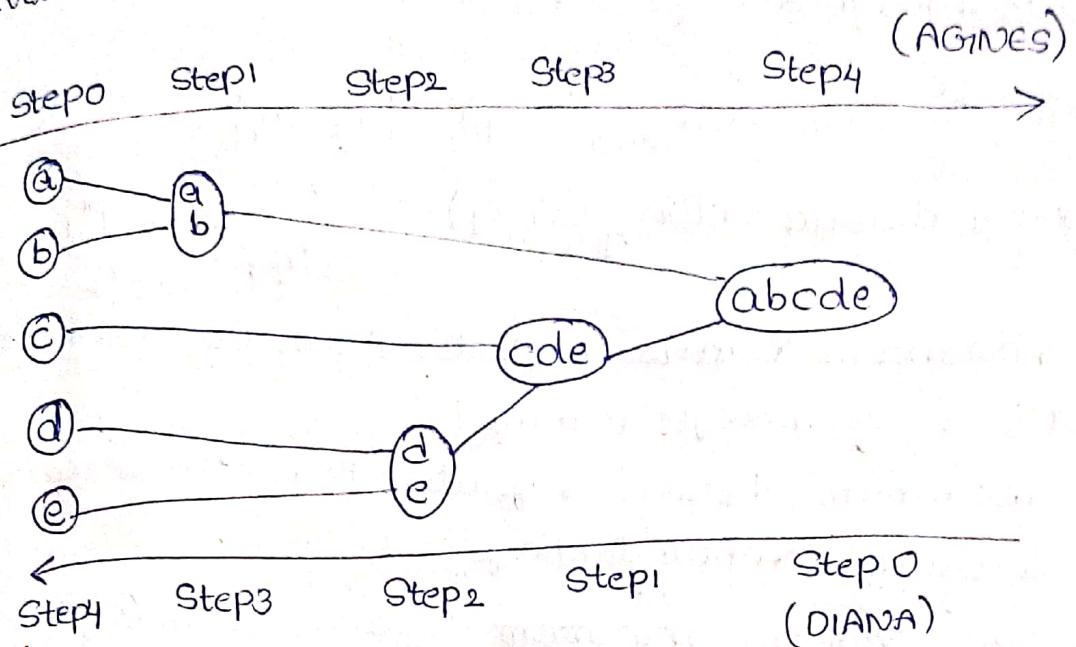
Agglomerative hierarchical clustering method :-

This method uses bottom up strategy Initially, each object having own cluster and each object combined with other object to form a new cluster until all the objects are in a single cluster. Or certain termination conditions are satisfied. The single cluster becomes the hierarchy root

Divisive hierarchical clustering Method :-

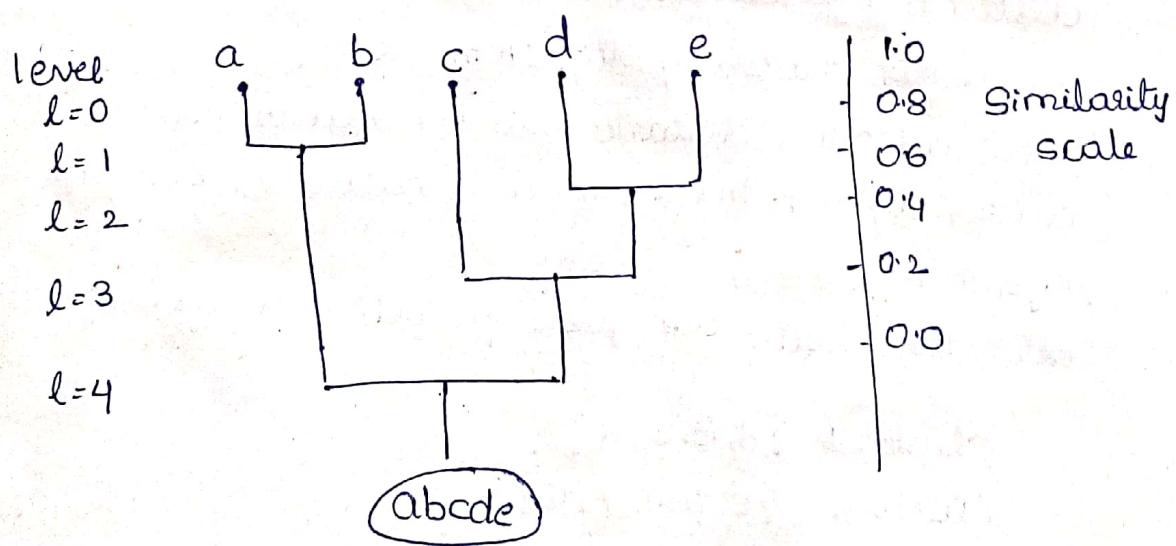
This method uses top down strategy it starts by placing all objects in one cluster which is a hierarchy root. It then divides the root cluster into several smaller sub clusters and recursively partitions those clusters into smaller ones. The partitioning process continues until each cluster at lowest level is coherent enough either containing only one object.

The Agglomerative hierarchical clustering also called as AGNES and divisive hierarchical clustering also called as divisive analysis Diana.



Dendrogram:-

A Tree structure called a dendrogram is commonly used to represent the process of hierarchical clustering. It shows how objects are grouped together on partitioned step by step.



Dendrogram Representation

Q) Distance measures in Alg methods :-

→ Nearest Neighbour

$$\text{minimum distance } \text{dist}_{\min}^{(C_i, C_j)} = \min_{p \in C_i, p' \in C_j} \{ |p - p'| \}$$

$$\text{Maximum Distance} \rightarrow \text{dist}_{\max}(c_i, c_j) = \max_{p \in c_i, p' \in c_j} |p - p'|$$

Farthest Neighbour

$$\text{Mean Distance: } \text{dist}_{\text{mean}}(c_i, c_j) = |\bar{m}_i - \bar{m}_j|$$

$$\text{Average distance: } \text{dist}_{\text{avg}}(c_i, c_j) = \frac{1}{n_i n_j} \sum_{p \in c_i, p' \in c_j} |p - p'|$$

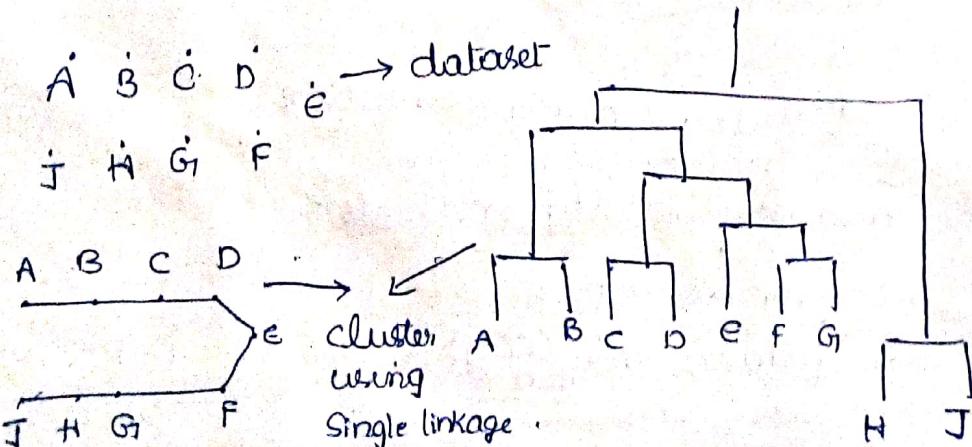
If minimum distance is greater than the threshold value is called single linkage.

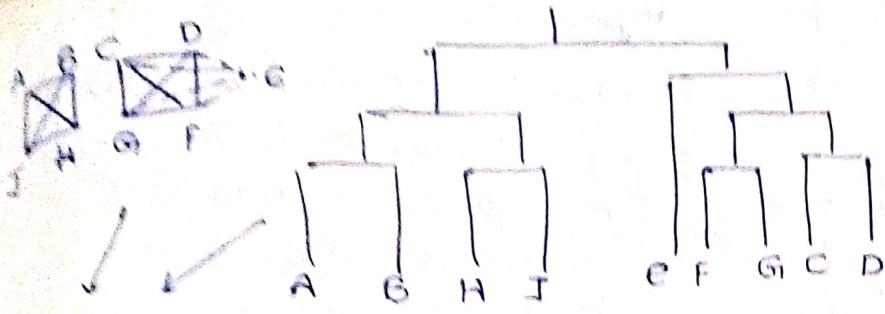
If maximum distance is greater than threshold value, it is called complete linkage.

- An algorithm uses minimum distance, $d_{\min}(c_i, c_j)$ to measure the distance b/w clusters, it is sometimes called as nearest neighbour clustering algorithm.
- When the min distance b/w nearest clusters exceeds a user defined threshold is called single linkage Algo.
- When an algorithm uses $d_{\max}(c_i, c_j)$ to measure the distance b/w clusters it is sometimes called as farthest neighbour clustering algorithm.
- When the max distance b/w the nearest clusters exceeds a user defined threshold it is called as complete linkage algorithm.
- Here the clusters are given for below diagrams.

Cluster 1 {A, B, I, H}

Cluster 2 {C, D, G, F, E}





cluster
using complete
linkage.

partitional based clustering :-

k-Means clustering :-

$K \rightarrow$ No. of clusters, means \rightarrow mean

$x = \{2, 4, 10, 12, 3, 20, 30, 11, 25\}$ apply K-means algo and
grp the data into two clusters

$$K=2 \quad m_1=7 \quad m_2=10$$

$$K_1 = \{2, 3\}$$

$$K_2 = \{4, 10, 12, 20, 30, 11, 25\}$$

$$m_1 = 2.5 \quad m_2 = 10$$

$$K_1 = \{2, 3, 4\} \quad K_2 = \{10, 12, 20, 30, 11, 25\}$$

	m_1	m_2	K_1	K_2
	3	18	{2, 3, 4, 10}	{12, 20, 30, 11, 25}
	4.75	19.6	{2, 3, 4, 10, 11, 12}	{20, 30, 25}
	7	25	{2, 3, 4, 10, 11, 12}	{20, 30, 25}

k-medoids algorithm (or) PAM Algorithm :-

PAM (partitioning Around medoids)

medoid is nothing but a data point which contains the most similarity with other datapoints in cluster

while making the data into clusters a medoid is a compulsory datapoint in a cluster

Ex:- Document	x	y	Algo :-
d ₁	2	6	1. Select k no of medoids randomly
d ₂	3	4	2. we calculate the distance from
d ₃	3	8	medoids to other data points
d ₄	4	7	3. calculate the total cost in making
d ₅	6	2	clusters
d ₆	6	4	4. If previous total cost is less than
d ₇	7	3	the present one while making the
d ₈	7	7	clusters stop the procedure.

Sol:- K=2 (3,4) (7,4)

Manhattan distance method $|x_1 - y_1| + |x_2 - y_2|$

	(3,4)	(7,4)	
d ₁	$ 2-3 + 6-4 = 3$	$ 2-7 + 6-4 = 7$	
d ₃	4	8	cluster 1 = {d ₁ , d ₃ , d ₄ }
d ₄	4	6	cluster 2 = {d ₅ , d ₆ , d ₇ }
d ₅	5	3	Total cost = 16
d ₆	3	1	we will choose some other
d ₇	5	1	medoids (3,4) (7,3)

Iteration 2:

	(3,4)	(7,3)	
d ₁	3	8	cluster 1 = {d ₁ , d ₃ , d ₄ }
d ₃	4	9	cluster 2 = {d ₅ , d ₆ , d ₈ }
d ₄	4	7	Total cost = 16
d ₅	5	2	
d ₆	3	2	
d ₈	4	1	

Hierarchical:-

By using following Techniques we develop hierarchical for in distance measure.

Clustering
using single link.

	x	y
P ₁	0.40	0.53
P ₂	0.22	0.38
P ₃	0.35	0.32
P ₄	0.26	0.19
P ₅	0.08	0.41
P ₆	0.45	0.30

1. findout the weighted distance matrix for the given data by using Euclidean distance
2. Find out the min value among all the distances b/w two data points & make them cluster.
3. continue the procedure until all the datapoints merge into single data cluster.

Step: P₁ P₂ P₃ P₄ P₅ P₆

P ₁	0				
P ₂	0.234	0			
P ₃	0.22	0.15	0		
P ₄	0.37	0.20	0.15	0	
P ₅	0.34	0.14	0.28	0.29	0
P ₆	0.23	0.25	0.11	0.22	0.39

cluster small.

P₃ P₆
cluster

P₁ P₂ {P₃ P₆} P₄ P₅

P ₁	0			
P ₂	0.234	0		
{P ₃ , P ₆ }	0.22	0.15	0	
P ₄	0.37	0.20	0.15	0
P ₅	0.34	0.14	0.28	0.29

For {P₃, P₆} $\rightarrow \min\{dist\{P_3, P_6\}\}$

{P₆, P₁} }

{0.22, 0.23} = 0.22

$\min\{(P_3, P_1), (P_6, P_2)\} = 0.15$

P₁ P₂P₅ P₃P₆ P₄

P₁

P₂P₅

P₃P₆

P₄

process continues until we get
single cluster

Nearest Neighbour clustering Method:-

It is an example for the partitioned based clustering Method. By using single link techniques, we calculate the clusters for given data sets.

In this Techniques, items are iteratively merged into the existing clusters which are closest.

In this Algorithm, a threshold value t is used to determine, If items will be added to existing clusters/not. If the distance b/w the two items is less than the threshold value, we club into single cluster.

Ex:- Threshold $t = 2$

data - $\{1, 2, 3, 4, 5, 7\}$

$$K_1 = \{ \underbrace{1, 2, 3} \}_{\substack{\text{Taken} \\ \text{value}}} \quad 2-1 = 1 < t$$

$$K_2 = \{ 4, 5 \} \quad 3-1 = 2 = t$$

$$K_3 = \{ 7 \}$$

Density Based clustering:-

clustering based on density (local cluster criteria)
such as density connected points are based on explicitly constructed density function.

Advantages -

- Discover clusters of different shapes
- Handle noise
- It uses single scan
- It requires Density parameters

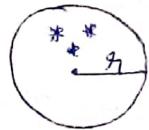
Ex:- DB Scan

DBSCAN :-

Density Based Spatial clustering of applications with noise

density : No. of points within a specified radius(ϵ) ~~area~~

- ϵ SP or ϵ or ϵ

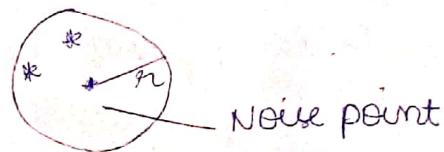
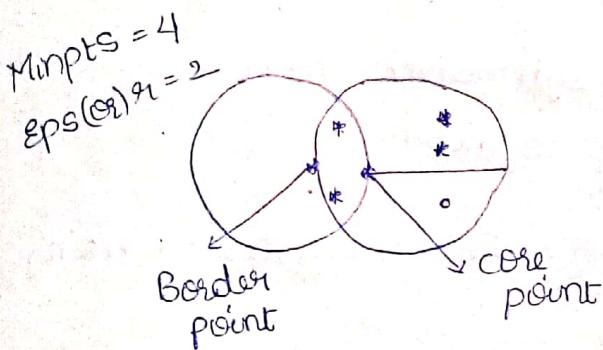


core point - A point is a core point, if it has more than specified no. of points ($minpts$) within ϵ

These are the points that are interior of a cluster

Border point - A border point has less than $minpts$ within ϵ but is in the neighbourhood of a core point

Noise point :- It is not a core point nor a border point



Algorithm :-

- for each $a \in D$ do
 - if a is not yet classified then
 - if a is a core point then
 - collect all objects density reachable from a
 - and assign them to a new cluster
 - else
 - assign to Noise

Density - Reachability :-

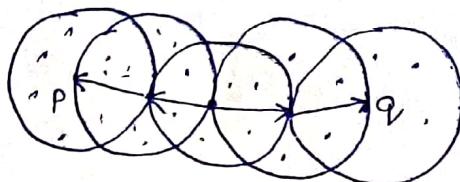
- A point q is directly density reachable from a point p if p is a core point and q is p 's ϵ -neighbourhood

Ex:-



Density connectivity :-

A pair of points p & q , are density-connected, If they are commonly density-reachable from a point x .



$$\text{minpts} = 7$$

Density-connectivity - Symmetric

Density-Reachability - Asymmetric

BIRCH - Balanced Iterative Reducing & clustering using hierarchies

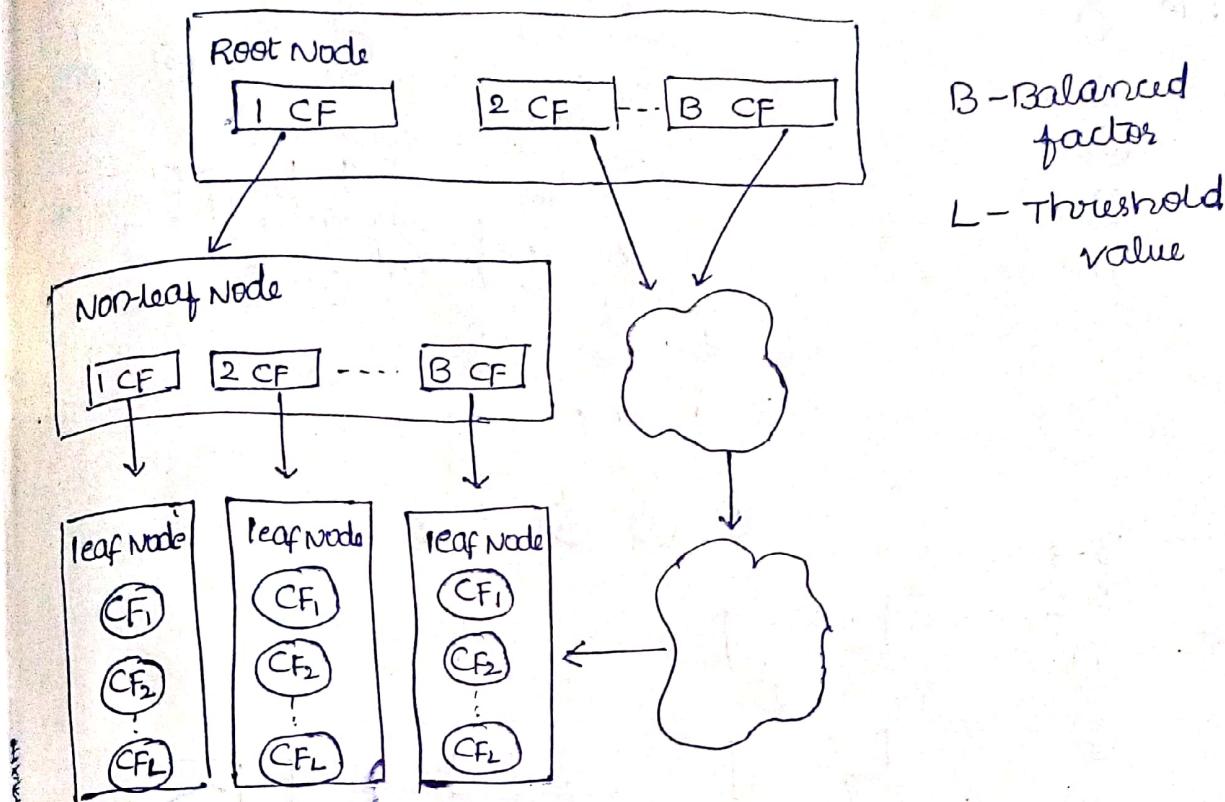
- It is the scalable clustering method
- Designed for very large datasets
- It requires only one scan of data if necessary
- It is based on the notation of CF (clustering feature) of a CF tree
 - A CF tree is a height balanced Tree that stores the clustering features for a hierarchical clustering
 - cluster of data points is represented by triple nos (N, LS, SS)

N = No. of items in sub-cluster

LS = Linear sum of the points

SS = Sum of the squares of the points

CF-Tree



B - Balanced factor
L - Threshold value

Height Balanced Tree with two parameters

1. Balancing factor (B) which indicates almost the pre-entries in non-leaf Node
2. (Each Leaf Node has almost) Threshold (L) which indicates each leaf Node has almost L entries (L CF entries)

The main advantage of BIRCH Algo is to handle memory restrictions

- It finds a good clustering with a single scan & improve quality.

Disadvantage - handles only numerical data

Applications :-

- pixel classification in Images
- Image compression
- Speech Enhancement

$$CF = (N, LS, SS)$$

$$LS: \sum_{i=1}^N x$$

$$SS: \sum_{i=1}^N x^2$$

$$\begin{aligned} Ex: (3, 4)(2, 6)(4, 5)(4, 7)(3, 8) \\ N=5 \end{aligned}$$

$$\begin{aligned} LS = 3+2+4+4+3 = 16 \\ 4+6+5+7+8 = 30 \end{aligned}$$

$$\begin{aligned} SS = 3^2+2^2+4^2+4^2+3^2 = 54 \\ 4^2+6^2+5^2+7^2+8^2 = 190 \end{aligned}$$

$$(5, (16, 30)(54, 190))$$

Ex:- $B=3$

$T=1$

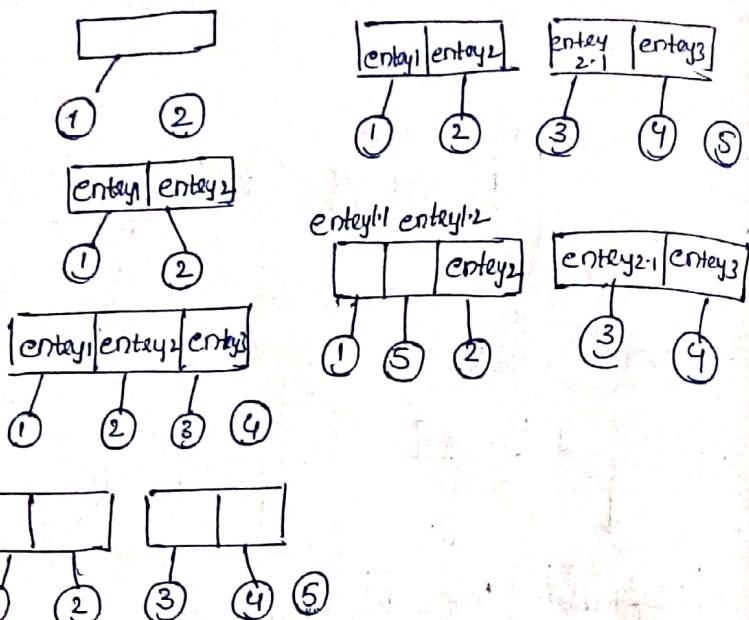
data ①

②

③

④

⑤



$B=3$

$T=2$

Ex:-

①

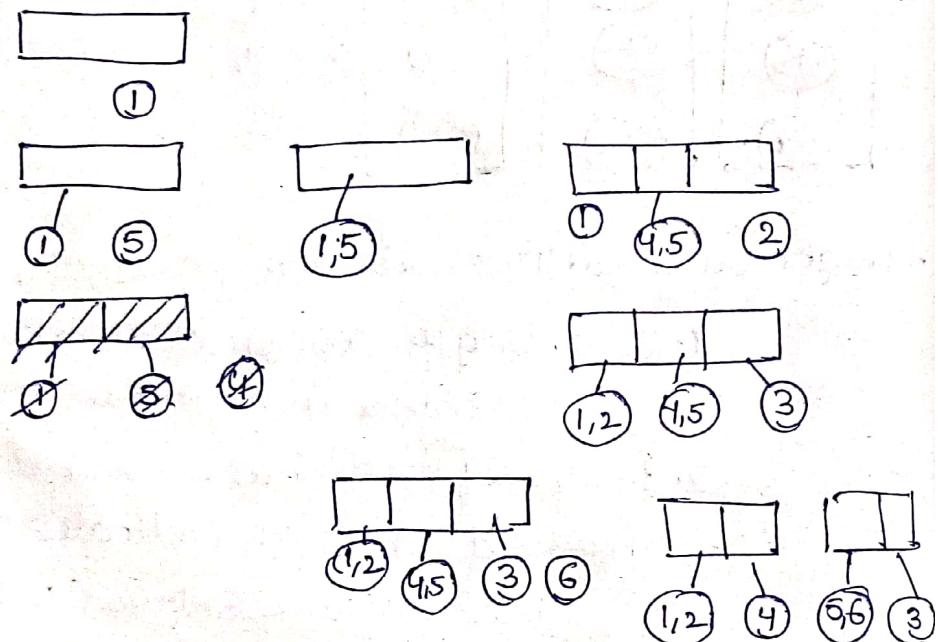
⑤

④

②

③

⑥



Similarity & Distance Measures :-

similarity indicates a tuple within one cluster is more similar within that cluster than it is similar to outside of clusters. by using distance measures, we can divide the data into the different clusters

$$\text{dist}(t_i, t_j) \rightarrow t_i \in K_i, t_j \in K_j, t_i \notin K_j, t_j \notin K_i$$

Techniques to find Distances :-

$$1. \text{Centroid } C_m = \frac{\sum_{i=1}^N (t_m)}{N}$$

$$2. \text{ radius } R_m = \sqrt{\frac{\sum_{i=1}^N (tm_i - m)^2}{N}}$$

$$3. \text{ diameter } D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (tm_i - tm_j)^2}{(N)(N-1)}}$$

- The Radius is the square root of the Average mean square distance from any point to clusters to the centroid.

- Diameter is the square root of the Average mean square distance b/w all pairs of points in the clusters.

The Drawback of these methods is Time consuming process and costly operation. So replace these methods by following

techniques - Single link

- complete link
- Average link
- centroid
- Medoid

CURE (Clustering with Representatives) :-

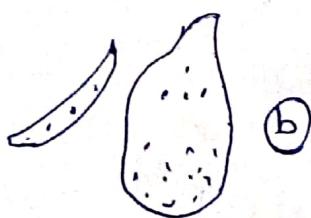
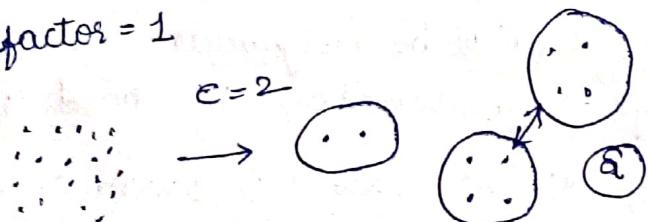
It is an example for both partitional & hierarchical

cluster.

A constant no. of points c is chosen from each cluster and shrink towards the each clusters, If shrinking

factor = 1

Ex:-



Clustering with categorical data :-

{ water, book, sun, sand, swim, read }

doc1 = { book } \approx [1, 0, 0, 0, 0, 0]

doc2 = { water, sun, sand, swim } = [0, 1, 1, 1, 1, 0]

doc3 = { water, sun, swim, read } = [0, 1, 1, 0, 1, 1]

doc4 = { read, sand } = [0, 0, 0, 1, 0, 1]

By using Euclidean distance

	1	2	3	4	
1	0				[2,3]
2	2.24	0			4
3	2.24	(1.41)	0		
4	1.73	2	2	0	

- But there is no similarity b/w doc1 & 4, so, we are going to ROCK algorithm.

ROCK Algo :- (Robust clustering with links)

- This Algo works for both categorical & boolean data.
- It is identifying the no. of links b/w the nodes (or) docs.
- If the no. of links greater than the threshold value, we can combine the data into single cluster.
- The pair of items are said to be neighbours, if they exceeds the links b/w objects greater than threshold value.
- Instead of Euclidean distance, we use Jaccard Quotient to find out the similarity b/w two documents.

$$\text{Sim}(t_i, t_j) = \frac{|t_i \cap t_j|}{|t_i \cup t_j|}$$

	1	2	3	4
1	1	0	0	0
2	0	1	0.6	0.2
3	0	0.6	1	0.2
4	0	0.2	0.2	1

$$2,3 \rightarrow \frac{3}{5} = 0.6$$

$$2,4 \rightarrow \frac{1}{5} = 0.2$$

$(Q, 4) (Q, 3) (3, 4)$ for threshold $t = 0.2$

Bond energy clustering Algo:-

- It is used in the database design to determine how to group the data and how to place physically data on to the disk.

- with BEA, Find out the affinity (bond) b/w the db attributes based on common usage.

- The basic objective of BEA Algo is Attributes that are used together form a cluster and should be stored together resulting cluster called as vertical fragments.

$$Aff = \sum_{i=1}^n (lbond(A_i, A_{i-1}) + bond(A_i, A_{i+1}))$$

clustering with Genetic Algo:-

- To apply this convert the given data into binary format

Ex:- $\{A, B, C, D\}$

\Downarrow \Downarrow
 $\{A, D\}$ $\{B, C\}$

\Downarrow \Downarrow
 $\{1, 0, 0, 1\}$ $\{0, 1, 1, 0\}$

Algorithm:-

Step 1: Randomly Create an initial solution.

Repeat:

use crossover to create a new step

until terminator criteria is met

$\text{ext} = \{A, B, C, D, E, F, G, H\}$

\Downarrow , 3 clusters

$\{A, C, E\} \quad \{B, F\} \quad \{D, G, H\}$

Cross-over :-

$\{A, B, C, D, E, F, G, H\}$

3-clusters

$\{A, C, E\} \quad \{B, F\} \quad \{D, G, H\}$

$\Downarrow \quad \Downarrow \quad \Downarrow$

$1010|1000 \quad 01000100 \quad 0001|0011$

bit-map representation

$\{A, B, C, D\}$

$\{A, D\} \quad \{B, C\}$

$\{1001\} \quad \{0110\}$

① & ③ are clusters apply crossover

$00011000 \quad 01000100 \quad 10100011$

$\{D, E\}$

$\{A, C, G, H\}$

Clustering with neural Networks :-

The neural networks that used unsupervised learning attempt to find the features in data that characterize the design O/P.

The neural networks search for the clusters with require links to fit the given data.

There are two types of neural networks in unsupervised learning - Non-competitive
- competitive

- Non competitive or Hebbian learning :-

The weight b/w two nodes is if change is proportional to the o/p values.

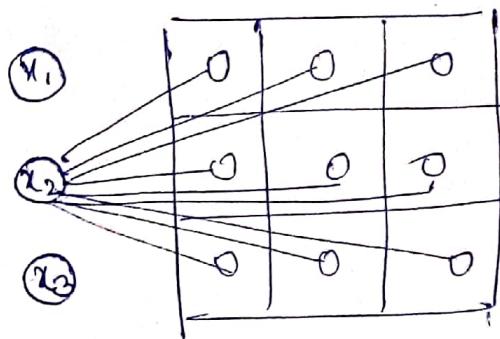
$$\Delta w_{ji} = \eta y_j y_i$$

These types of n/w also called as self-organising neural networks

To construct self-organising neural n/w, we use self-organising feature maps.

This algorithm works like the behavior of a node "should impact only those nodes and edges which are nearer to that node".

The best ex of Self-organising neural n/w is Kohonen Self-organising map



$$\text{sin}(x, i) = \sum_{j=1}^n x_j w_{ji}$$

$$\Delta w_{kj} = \begin{cases} c(x_k - w_k) & \text{if } j \in N_i \\ 0 & \text{otherwise} \end{cases}$$

In Kohonen Self Organising map there are two layers.

1. Input layer
2. competitive layer

- In competitive layer the nodes are arranged in 2D format every node is associated with activation function.
- every layer in the competitive one produces an output and which layer output is equal to the expected output make all the nodes in the layer into single cluster