

11/9/19

Cluster Analysis

For ex: At a bank,
 customers - 100 (Age, Income, Area, Gender)
 Cust. manager - 4 ← Input para. consist of
 4 clusters as clusters → the 4 attr.
 of cust. mentioned

- Clustering is the process of grouping the data into classes or clusters so that objects within cluster have high similarity in comparison to one another but, varying dissimilarity with objects of other clusters.
- Clustering gives dense or sparse representation of the data.
- Generally, in statistical distance-based clustering algo. are used.
- Classification classifies the given data into pre-defined labels only. It is supervised learning.
- Whereas, clustering is unsupervised learning.

Requirements of clustering in DM

- 1) Scalability - The algo. Should be able to work on any no. of data / tuples.
- 2) Ability to deal with different type of attributes - Multiple attr. Should be

predicted simultaneously.

- 3) Identifying clusters with arbitrary shape
- 4) Minimal / less knowledge for input parameters
 - cluster should be identified by using a min. no. of attr. or knowledge base
- 5) Deal with noisy data
 - Preprocessing should be performed in cases where noisy data is present.
- 6) Insensitivity towards the ordering of input parameters
 - whether data is processed from 1 to 100 or 100 to 1, the output should be same.
 even if some data is added, it should be processed in a similar way as its predecessors.
- 7) High Dimensionality
 - should be able to handle data from multiple tables or data in same table with multiple attr.
- 8) Constraint-based clustering
- 9) Interpretability & Usability
 - output of clustering should be

easy to interpret and easy to use

Type of data used in cluster analysis:

1) Data matrix

Represented as 'Object - by - variable'

Person Age, Income
n P

∴ $(n \times p)$ matrix would be generated.

$$\begin{bmatrix} x_{11} & \dots & x_{1p} \\ x_{21} & \dots & x_{2p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$

2) Dissimilarity matrix

- Represented as 'Object - by - Object'
- It shows the dissimilarity between two objects

Same as $d(2,1)$

	1	2	3	\dots	n
1	D				
2	$d(2,1)$	D			
3	$d(3,1)$	$d(3,2)$			
n	$d(n,1)$	\dots			D

Major clustering methods

1) Partitioning method

2) Hierarchical method

3) Density-based method

4) Grid-based method

5) Model-based method

Partitioning Method.

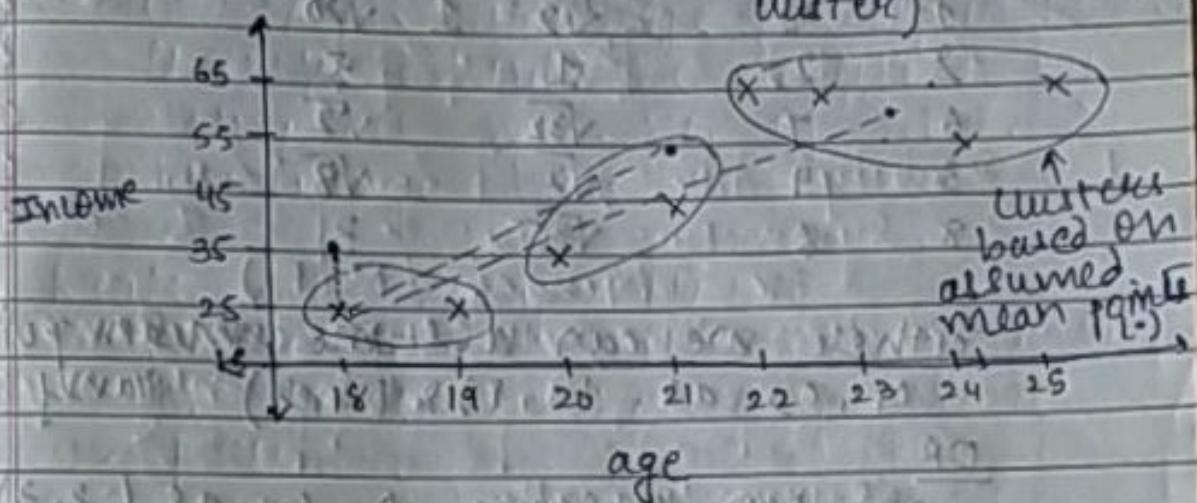
- Distance-based algo. used.

i) K-mean

ii) K-medoids

// K = no. of clusters.

(mean = mean of particular cluster)



Assume k=3

- Mark 3 (since k=3) random points on the graph
- Measure distance of every point (x) from the three random points (•) (Use Euclidean dist. formula)
- Cluster the points having roughly the same avg. dist. from the random point (•)
- These steps are iterated for various 3-random points until a min. dist. from the mean of cluster is obtained.
(Typically, iterated \approx times)

<u>Ex:</u>	(2, 3)	(1, 4)	(3, 4)	$k=2$
	(5, 6)	(2, 2)	(8, 6)	
	(8, 7)	(6, 7)		

\Rightarrow	<u>x</u>	<u>y</u>	$C_1(2,3)$	$C_2(5,6)$	cluster
	2	3	0	$9\sqrt{2}$	g
	5	6	$9\sqrt{2}$	0	C_2
	8	7	$\sqrt{52}$	$10\sqrt{10}$	C_2
	1	4	$\sqrt{2}$	$\sqrt{20}$	C_1
	2	2	1	5	C_1
	6	7	$\sqrt{32}$	$\sqrt{2}$	C_2
	3	4	$\sqrt{2}$	$\sqrt{8}$	C_1
	8	6	$\sqrt{45}$	3	C_2

either we can assume mean points
 c_1, c_2 as $(1,1)$ & $(3,4)$ (any random)

OR

we can assume c_1, c_2 as $(2,3)$ & $(5,6)$

new c_1

c_2

$(2,3)$

$(5,6)$

$(1,4)$

$(8,7)$

$(2,2)$

$(6,7)$

$(3,4)$

$(8,6)$

$$c_1 = \left(\frac{2+1+2+3}{4}, \frac{3+4+2+4}{4} \right) \text{ new } c_1 = \left(\frac{27}{4}, \frac{26}{4} \right)$$

$$= (2, 8.25) \quad \text{new } c_2 = (6.75, 6.5)$$

Iteration 2:

<u>x</u>	<u>y</u>	$C_1(2,3.25)$	$C_2(6.75,6.5)$	cluster
2	3	0.25	6.9	C_1
5	6			
8	7			
1	4			
2	2			
6	7			

3 4
8 6

13/9/18

- Types of attributes possible in a cell
- Nominal attr. (related to names/categories)
 - Binary attr.
 - Ordinal
 - Numerical
 - Mixed

Proximity measure - How the obj. are alike or unlike when compared to one another.

(i) \Rightarrow Let assume two objects, i & j
 $d(i, j) = \frac{P-m}{P}$ $s(i, j) = 1 - d(i, j)$
[similarity measure]

P = Total no. of attr.

m = Total no. of matched attr.

object ID color grade

1 Red A

2 Yellow B

3 Red A

4 Blue C

$$d(1, 2) = \frac{2-0}{2} = 1 \quad s(1, 2) = 0,$$

$$d(1, 3) = \frac{2-2}{2} = 0$$

$$d(1, 4) = \frac{2-0}{2} = 1$$

$$\delta(2,3) = \frac{2-0}{2} = 1$$

$$\delta(2,4) = \frac{2-0}{2} = 1$$

$$\delta(3,4) = \frac{2-0}{2} = 1$$

(ii) Binary attr.

- symmetric attr.

Ex - Gender (M, F)

Both represent same attr.

& can take only three values

- asymmetric attr.

Ex - +ve/-ve (+ve, -ve)

object i

1 0

Object j

1 q t

0 s t

Symmetric binary

$$\delta(i,j) = \frac{q+s}{q+s+t+t}$$

Asymmetric binary

$$\delta(i,j) = \frac{s+t}{q+s+t+t}$$

Jaccard coefficient

Name	Fever	Cough	T1	T2	T3	T4
Jack	Y	N	P	N	N	N
Mary	Y	N	P	N	P	N
Jim	Y	Y	N	N	N	N

All are asymmetric attr.

Jack

$$\begin{array}{cc} 2 & 1 \\ 0 & 3 \end{array}$$

$$d(\text{Jack, Mary}) = \frac{1+0}{2+1+0} = \frac{1}{3} = 33.3\%$$

Jack

$$\begin{array}{cc} 1 & 1 \\ 1 & 3 \end{array}$$

$$d(\text{Jack, Jim}) = \frac{1+1}{1+1+1} = \frac{2}{3} = 66.6\%$$

Mary

$$\begin{array}{cc} 1 & 1 \\ 2 & 2 \end{array}$$

$$d(\text{Mary, Jim}) = \frac{1+2}{1+1+2} = \frac{3}{4} = 75\%$$

(iii) Ordinal attr.

Order is given as attr. value
ex - low, medium, high or 1st, 2nd, 3rd

ID	Grade	ID	Grade
1	High	1	1
2	Low	2	0.5
3	High	3	0
4	Med	4	0.5

$$R_{if} = \{1, 2, 3, \dots, M_f\} ; M_f = \text{Total stages}$$

(Here, 3)

$$R_{if} = \{1, 2, 3\}$$

High ↑ Med ↑ Low ↑

$$z_{if} = \frac{R_{if} - 1}{M_f - 1}$$

$$\therefore \text{for high, } z_{if} = \frac{1 - 1}{3 - 1} = 0.5$$

Page No.	
Date	/ /

Using Manhattan dist.

$$d(A, B) = |1| + |-2| + |-10| + |5| \\ = 1 + 2 + 10 + 5 \\ = 18$$

16/9/19

a) Interval Scaled

- linear measurement

Temp

30°

60°

0°

→ is not half of 60°. Just
represents a range

b) Ratio Scaled

Ex:

Income

25K

50K

75K

"Here, 50K is considered
as double of 25K.
(i.e. ratios are
involved)

converting interval scaled values to
standardized values.

1) Mean absolute deviation

$$S_f = \frac{1}{n} (|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

var. of
attr.

m_f = mean value of f

$$= \frac{1}{n} (x_{1f} + x_{2f} + \dots + x_{nf})$$

2) Standardized measurement.

$$z_{if} = \frac{x_{if} - m_f}{S_f}$$

(calculated for
each value of
an attr.)

Generalized
Combining eq. of Euclidean + Manhattan:

$$d(i, j) = [(|x_{i1} - x_{j1}|)^q + (|x_{i2} - x_{j2}|)^q + \dots + (|x_{ip} - x_{jp}|)^q]^{1/q}$$

$q=1$: Manhattan dist.

$q=2$: Euclidean dist.

Ex: calculate dissimilarity matrix for:

A (7, 8, 3)

B (9, 7, 5)

C (7, 6, 1)

\Rightarrow Total 6 matrices would be generated

$$d(A, B, C) = \begin{matrix} A & \begin{bmatrix} 0 & 3 & \sqrt{8} \\ B & 0 & \sqrt{29} \\ C & & 0 \end{bmatrix} \end{matrix}$$

$$\rightarrow d(i, j) \geq 0$$

$$d(i, i) = 0$$

$$d(i, j) = d(j, i)$$

$$d(i, j) \leq d(i, h) + d(h, j)$$

$$d(A, B) = \sqrt{2^2 + 1^2 + 2^2} = 3$$

$$d(B, C) = \sqrt{2^2 + 3^2 + 4^2} = \sqrt{29}$$

$$d(A, C) = \sqrt{0 + 2^2 + 2^2} = \sqrt{8}$$

$$d(A, A) = d(B, B) = d(C, C) = 0$$

Ratio-scaled (3 methods)

- i) Treat the ratio-scaled variable like the interval-scaled variable.
- ii) Apply logarithm function for attr. value & then apply interval scaled technique.
- iii) Convert attr. values to continuous ordinal values & then apply (interval scaled) techniques.

(V) Mixed Variable

$$d_{ij}^{(f)} = \frac{\sum_{f=1}^P s_{if}^{(f)} \times d_{if}^{(f)}}{\sum_{f=1}^P s_{if}^{(f)}}$$

f: attr. value

where,

$s_{if}^{(f)} = 0$ if 1) x_{if} or x_{if} is missing
2) $x_{if} = x_{if} = 0$
& f is asym. binary

Otherwise,

$$s_{if}^{(f)} = 1$$

\Rightarrow 1) f is binary or nominal

$$d_{if}^{(f)} = 0, \text{ if } x_{if} = x_{if}$$

$$d_{if}^{(f)} = 1, \text{ otherwise}$$

2) f is interval scaled

$$d_{if}^{(f)} = |x_{if} - x_{if}|$$

maxn $x_{if} - \min x_{if}$

3) f is ordinal or ratio scaled

$$x_{ij} \& z_{ij} \in \frac{k_{ij}-1}{m_j-1}$$

Conditions for partitioning:

- i) Every object must belong to only one cluster.
- ii) Every cluster must have at least one object.

k-Medoids

The element having minimum dissimilarity with other ele. of the cluster.

	c_1	c_2	Cluster
A(8,7)	$\sqrt{2}$	$\sqrt{20}$	C_1
B(3,7)	$\sqrt{37}$	$\sqrt{5}$	C_2
C(4,9)	$\sqrt{34}$	4	C_2
D(19,6)	0	$\sqrt{25}$	C_1
E(8,5)	$\sqrt{2}$	4	C_1
F(5,3)	5	$\sqrt{5}$	C_2
G(7,3)	$\sqrt{13}$	$\sqrt{13}$	C_1, C_2
H(8,4)	$\sqrt{5}$	$\sqrt{17}$	C_1
I(7,5)	$\sqrt{5}$	3	C_1
J(4,5)	$\sqrt{26}$	0	C_2

$k=2$

$$K_1 = \{ A, E, G, H, I \}$$

$$K_2 = \{ B, C, F, J \}$$

$$\text{Cost.} = \sum \{ (DA), (DE), (DG), (DH), (DJ) \}$$

$$+ \sum \{ (JA), (JC), (JF) \}$$

Page No.	
Date	/ /

$$D_1 = (\sqrt{2} + \sqrt{2} + \sqrt{13} + \sqrt{5} + \sqrt{26}) + (\sqrt{5} + 4 + \sqrt{5}) \\ D_1 = 19.37$$

Repeat the steps by choosing a different c_1 or c_1 or both c_1 & c_2 .

If D_2 comes out to be more than D_1 , then it is considered as a bad decision and we stick with the earlier one.

18/9/19

PAM (Partitioning around medoid)

k clusters in object

$k=2$

Random Selection: c_1, c_2

Next, 100 new non-medoid points are taken as medoids and the dist from each point calculated. If it comes out to be less than the previous medoid, the new ones are considered & the process iterated.

(PAM)

Variations of k -medoid: k -median.

k -mode.

→ In presence of noise/outliers, k -medoid works better; else k -mean is mostly preferred.

Both the methods aren't sufficient of fulfilling when we have large amt. of data.

(2) CLARA (clustering Large Application) Sampling + PAM

The selection of sample is important

$$1-50 - 51-100$$

15 samples 5 samples // not good situation

Should be balanced

$$O(kl^2 + k(n-k)) \quad // \text{complexity of CLARA}$$

where, $s = \text{Sample size}$

$k = \text{no. of clusters}$

$n = \text{total no. of objects}$

Improving efficiency of CLARA : Use CLARANS.

Here, new (random) samples are generated for every iteration

[In CLARA, samples are fixed for every iteration].

Hierarchical clustering (method)

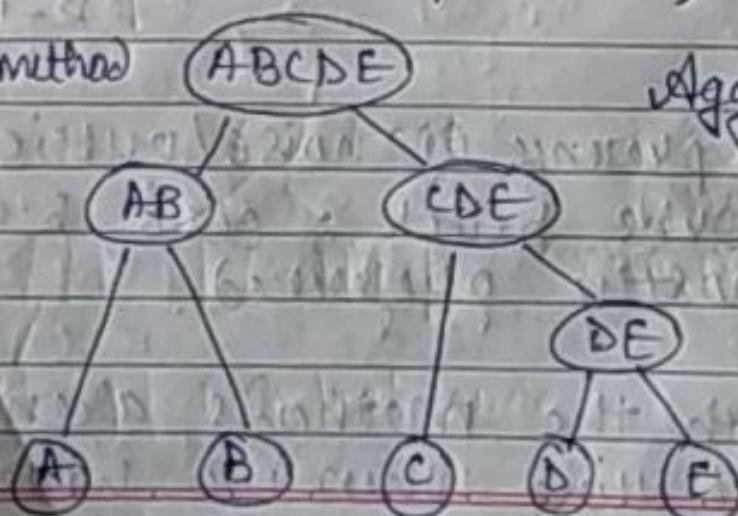
1) Agglomerative (Bottom-up)

2) Divisive (Top-down)

Divisive method

ABCDE

Agglomerative



(c) CLARA (clustering Large Application)
Sampling + PAM

The selection of samples is important
1-50 - 51 - 100

15 samples 5 samples // not good

solution should be balanced

$O(kl^2 + k(n-k))$ // complexity of CLARA

where, $S = \text{Sample size}$

$k = \text{no. of clusters}$

$n = \text{total no. of objects}$

Improving efficiency of CLARA : Use CLARANS.

Here, new (random) samples are generated for every iteration.

[In CLARA, samples are fixed for every iteration].

Hierarchical clustering (method)

1) Agglomerative (Bottom-up)

2) Divisive (Top-down)

Divisive method

ABCDE

AB

CDE

DE

A

B

C

D

E

Agglomerative



algo. for divisive : DIANA
 (Divisive Analysis)

algo. for Agg. : AGNES
 (Agglomerative Nesting)

No. of clusters is the termination point
 for both the methods.

Dist. b/w 2 clusters: $|P - P'|$ $P \in C_1$
 $P' \in C_2$

Four dist. to be considered:

i) Min. distance

$$c_i, c_j \\ \min(c_i, c_j) = \min(|P - P'|)$$

ii) Max. dist.

$$\max(c_i, c_j) = \max(|P - P'|)$$

iii) Mean dist.

$$\text{mean}(c_i, c_j) = [m_i - m_j] \quad \text{mean of } c_i$$

iv) Avg. dist.

$$\text{avg}(c_i, c_j) = \frac{1}{n_i n_j} \sum_{P \in c_i} \sum_{P' \in c_j} |P - P'|$$

Limitations:

1) Selecting splitting & merging point.
 (cannot be undone)

2) Scalability

(splitting & merging are costly operations)

Efficiency of hierarchical method can be improved by combining it with another method: BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies)

- Integrated Hierarchical Method

- i) clustering feature (CF)
- ii) clustering feature tree (CFT)

CF

- summarization triplets which include three factors:

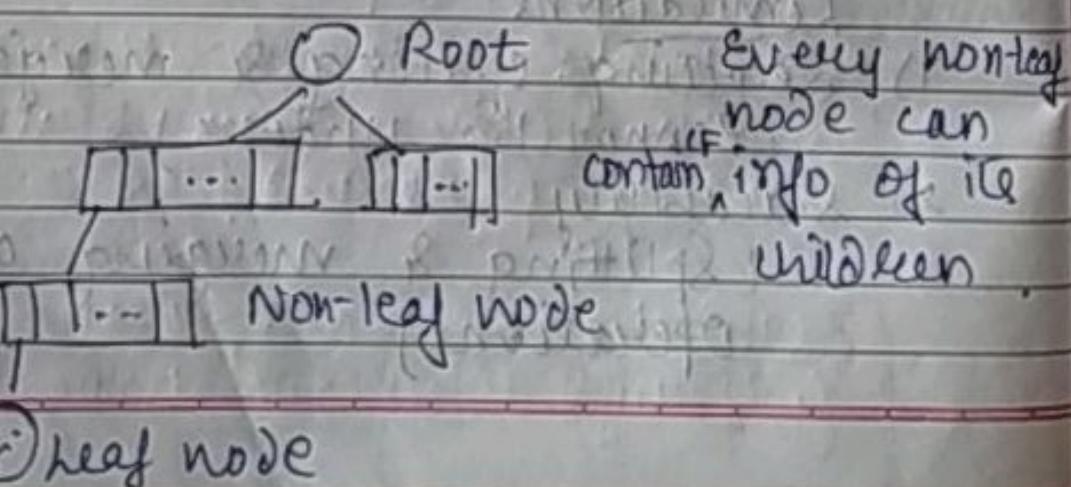
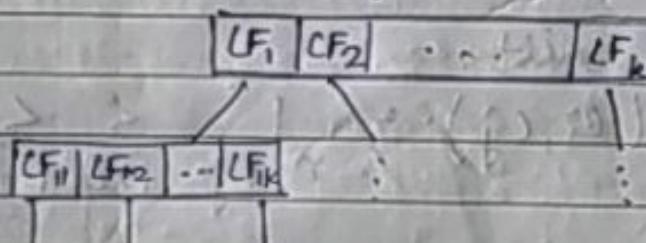
$$(N, \overrightarrow{IS}, SS)$$

NO. of linear
Obj. sum

sq. sum

$$\left. \begin{array}{l} A(3, 2) \\ B(5, 4) \end{array} \right\} CF = (2, (8, 6), (34, 20))$$

CF-tree (height-balanced tree)



Imp factors of BIRCH:

Branching Factor (B) : No. of children.
Threshold (D) Diameter of children.

⇒ BIRCH algo. phases:

second
Scan
can be
used to
remove
noise.

Phase-1 : Scan the database (dataset) & build the LF-tree. It is a dynamic process. (Data is added into the tree as it is read from the dataset). Complexity: $O(n)$

Phase-2 : Apply selected clustering algo. Data is added in the nearest cluster (at leaf node). When diameter increases, we can further split the cluster.

LURE of clustering using Representative: Centroid can be used as a representative for every cluster.

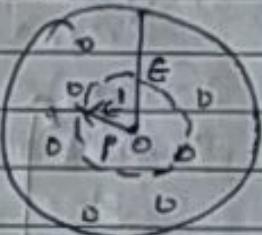
steps:

- 1) Draw random shape, S
- 2) Create partitions of the shape.
- 3) Cluster each partition (Find centroid of each cluster).
- 4) Eliminate outliers by random sampling of the clusters. (Entire cluster is eliminated as outlier)
- 5) Cluster the partial clusters.
- 6) Mark appropriate cluster with req. Obj (Shrink shape of cluster by finding mean centroid from centroids of sub cluster)

25/9/19

OPTICS

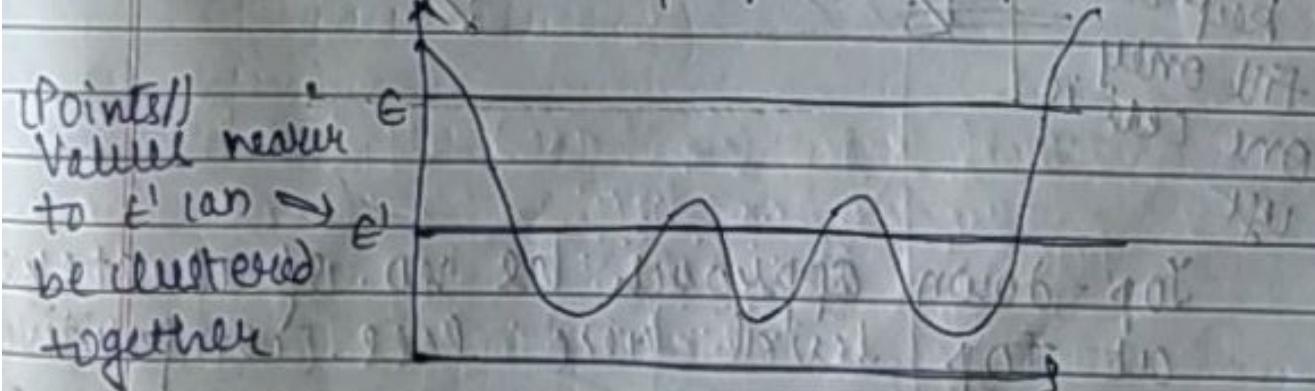
- 1) Core distance e'
 Min. dist. to make P
 the core point)



P = core-point
 min. dist.

- 2) Reachability distance
 (Max. value between (among)
 Euclidean dist. & core
 dist.)

After calculating for all points:



The ordering algo. on any attr.
 automatic/user-defined

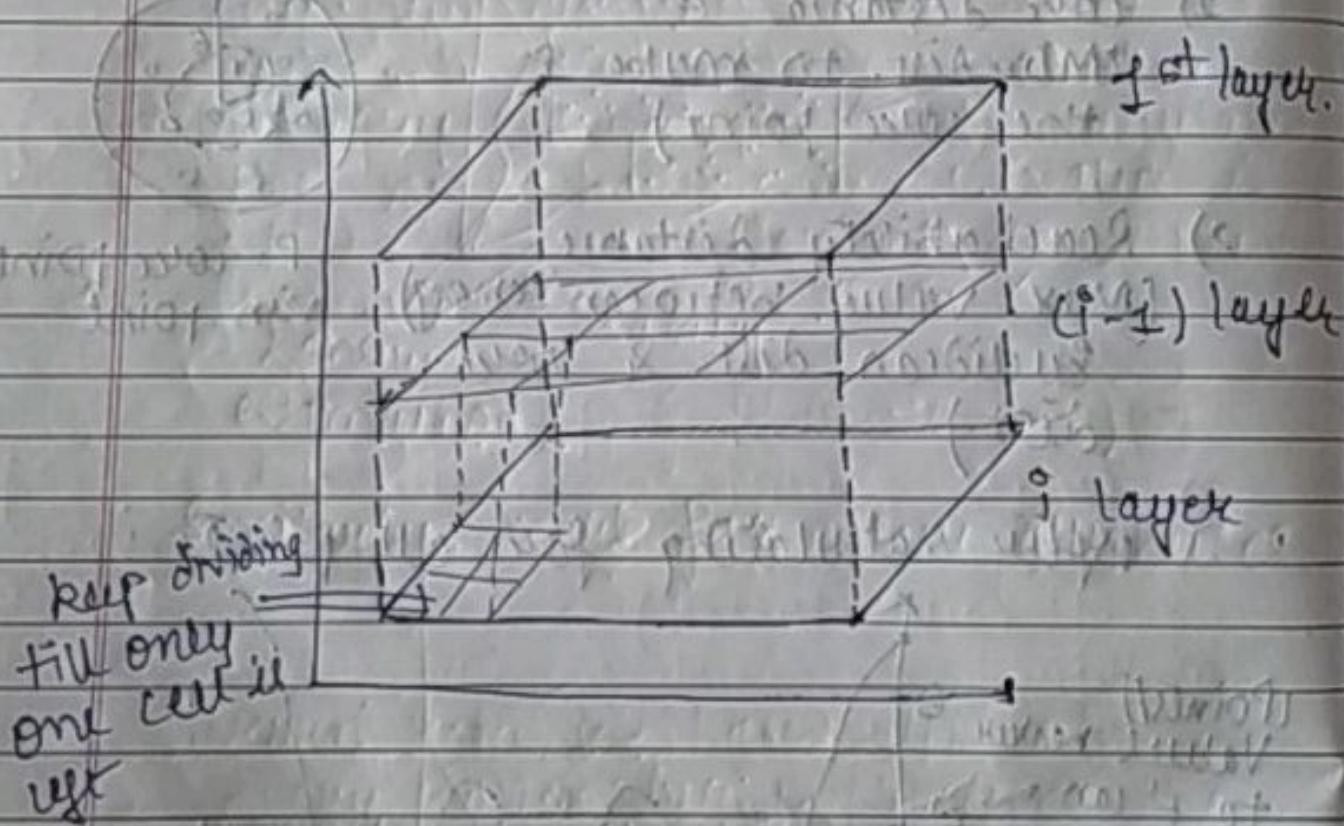
Any random shape can be generated when density-based clustering is used.

- 4) Grid-based clustering
Square shaped clusters can be generated.

Methode:

STING (Statistical Information Grid)
Combination of hierarchical and grid based clustering.

Make use of statistical data such as min., max. and mean. Also, SD.



Top-down approach as no. of clusters at top level are less (comparitively).

For each level, confidence interval $\frac{1}{k}$

~~calculated~~

eliminate irrelevant cells based on the confidence interval. Every level needs to store the statistical info. This increases the independence between levels.

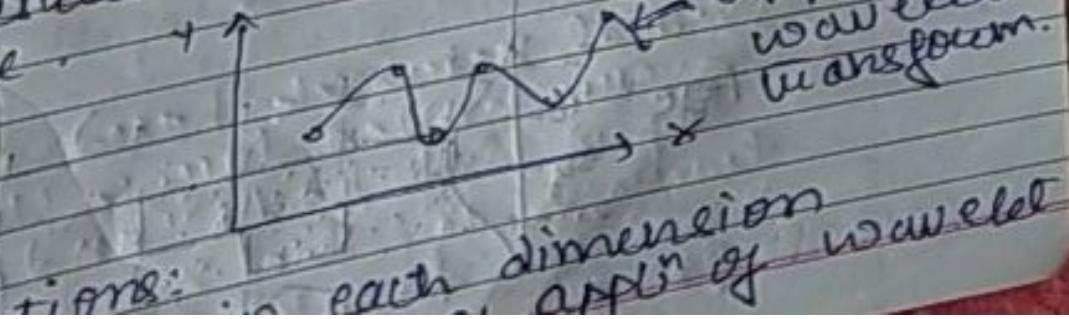
Wavelet transformation - Represent a wave as a sum of waves of various frequencies.

Clustering using wavelet tree.

all points in this part are considered as a cluster.

Noise

form a wave which includes all the points in the dataset. Then apply wavelet transform and compare the wave with the above figure.

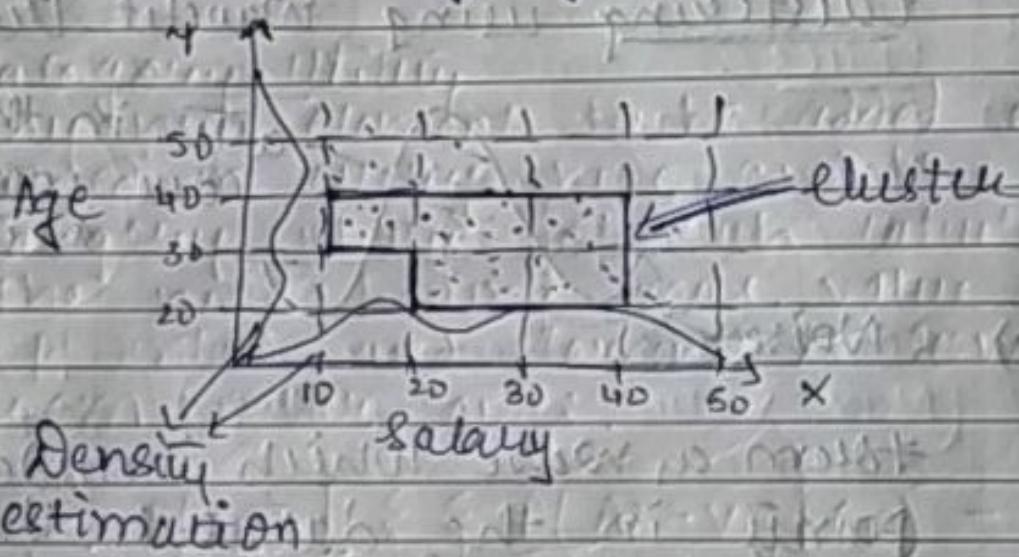


Arbitrary shape clusters can be generated.

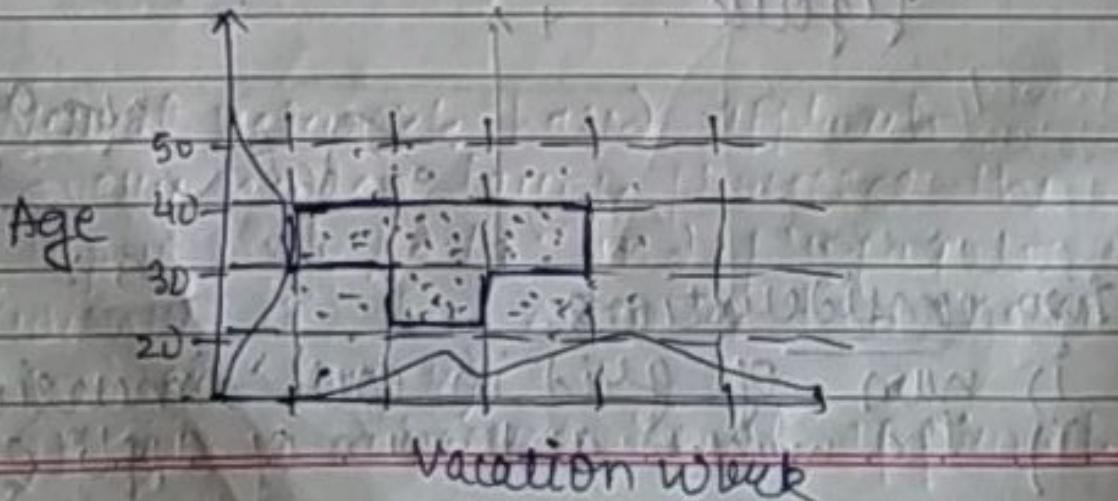
(Order insensitive algo. i.e. the order does not matter (as the points are plot in a grid & the order of plotting does not matter))

[combination of grid + density based clustering]

CLIQUE (clustering by-dimensional space)
used for high dimensional data
combination of grid + density based



Similarly, mark for other attributes



Final take intersection of both the clusters.

5) Model-based clustering

Mixture of unsupervised & supervised learning

(The other four methods all fall under the category of unsupervised learning)

It tries to fit given data into some mathematical model.

The output of first step (unsupervised learning) are clusters which are then used to build a model (supervised learning).

i) Statistical method / approach

- Uses conceptual learning.

Clusters \rightarrow concept of the clusters

extracting features of the class

Conceptual clustering:

I/P : Unlabelled objects (clusters)

O/P : Classification model based on characteristic / feature of obj.

COBWEB

$$P(A_i = v_{ij} | c_k)$$

\uparrow Attribute value

(attribute-value pair)

\leftarrow Concept class.

Classification, has two parameters:

i) concept

ii) Probabilistic description of that concept

Intra-class similarity:

$$P(A_i = v_{ij} | C_k)$$

Inter-class dissimilarity

$$P(C_k | A_i \neq v_{ij})$$

Animal

$$P(C_0) = 1$$

$$P(\text{scales} | C_0) = 0.25$$

Bird/Mammal

$$P(C_1) = 0.5$$

$$P(\text{hair} | C_1) = 0.5$$

fish

$$P(C_2) = 0.5$$

$$P(\text{scales} | C_2) = 1$$

Category utility

Any new element to be added is added in every node and the category utility for each calculated.

The node having max. value is kept intact and element is removed from other nodes.

$$\sum_{k=1}^n P(C_k) \left[\sum_j P(A_i = v_{ij} | C_k)^2 - \right.$$

$$\left. \sum_j P(A_i \neq v_{ij})^2 \right]$$

Neural network based approach

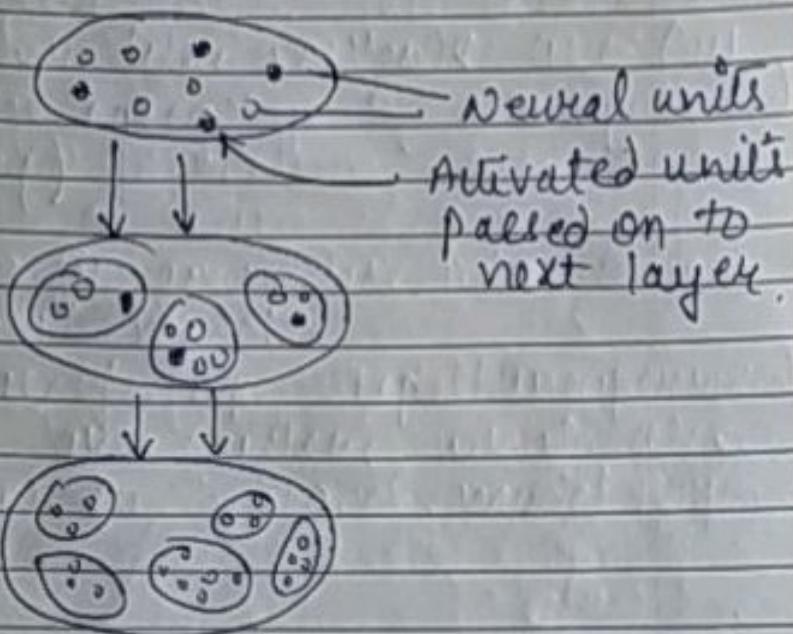
Every hidden layer processes the input info./data and propagates it further in the network (to next hidden layer).

It is a hierarchical approach

but, cannot be traversed bottom to top (i.e. from O/P layer to I/P layer)

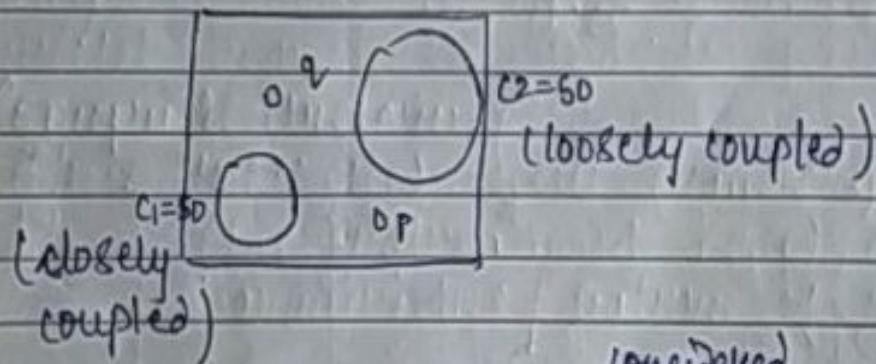
Types of methods: a) competitive
b) SOM, Self-organizing

Input



Outlier analysis

- a) Statistical methods
- b) Distance based detection



(one-sided)
universal threshold is used while using distance based. But, P may not be detected as an outlier as it is nearer to C2.