

# Cluster Analysis

(1)

2) Data Types in Cluster Analysis

3)

Categories of clustering methods

- Partitioning method  $\leftarrow$  K-mean  
K-means++  
CLARANS

- Hierarchical clustering

Agglomerative (bottom-up)  $\rightarrow$  ~~CURE~~ & Chameleon

Divisive (Top-down)  
BIRCH, Chameleon

- Density Based Methods

- DBSCAN

- OPTICS

- Grid-Based Methods

- STING

- CLIQUE

- Model Based Method

- Statistical approach

- Neural NW approach

- Outlier analysis

Applications -

1) Pattern Recognition

2) Spatial Data Analysis

3) Image Processing

4) Market Research

5) Document Classification

(2)

A good clustering method will produce high quality clusters with

- high intra-class similarity
- low inter-class -

⇒ Requirements of clustering in Data Mining

1) Scalability

2) Ability to deal with different types of attributes.

3) Ability to handle dynamic data

4) Discovery of clusters with arbitrary shape

5) Able to deal with noise & outliers

6) High dimensionality

7) Interpretability & usability

Cluster Analysis

K-Mean Clustering

- $A_1 \rightarrow (2, 10)$
- $A_2 \rightarrow (2, 5)$
- $A_3 \rightarrow (8, 4)$
- $A_4 \rightarrow (5, 8)$
- $A_5 \rightarrow (7, 5)$
- $A_6 \rightarrow (6, 4)$
- $A_7 \rightarrow (1, 2)$
- $A_8 \rightarrow (4, 9)$

Initial cluster centers are: (3)

- $A_1 (2, 10)$
- $A_4 (5, 8)$
- $A_7 (1, 2)$

Point	(2,10)	(5,8)	(1,2)	Cluster
(2,10)	0	5	9	1
(2,5)	5	6	4	3
(8,4)	12	7	9	2
(5,8)	5	0	10	2
(7,5)	10	5	9	2
(6,4)	10	5	7	3
(1,2)	9	10	0	3
(4,9)	3	2	10	2

Cluster 1  
(2,10)

Cluster 2  
(8,4)  
(5,8)  
(7,5)  
(6,4)  
(4,9)

Cluster 3  
(2,5)  
(1,2)



Center of the clusters:

(4)

1) Cluster 1  $\rightarrow$  (2, 10)

Cluster 2  $\rightarrow$  (6, 6)  $\left[ \frac{8+5+7+6+4}{5} \right], \left[ \frac{4+8+5+4+9}{5} \right]$

Cluster 3  $\rightarrow$  (1.5, 3.5)

(2, 10)

(6, 6)

(1.5, 3.5)

(2, 10)

(2, 5)

(8, 4)

(5, 6)

(7, 5)

(6, 4)

(1, 2)

(4, 9)

2<sup>nd</sup> epoch  $\rightarrow$

C<sub>1</sub> (3, 9.5) (A<sub>1</sub>, A<sub>8</sub>)

C<sub>2</sub> (6.5, 5.25) (A<sub>3</sub>, A<sub>4</sub>, A<sub>5</sub>, A<sub>6</sub>)

C<sub>3</sub> (1.5, 3.5) (A<sub>2</sub>, A<sub>7</sub>)

3<sup>rd</sup> epoch -

C<sub>1</sub> (3.66, 9) (A<sub>1</sub>, A<sub>4</sub>, A<sub>8</sub>)

C<sub>2</sub> (7, 4.33) (A<sub>3</sub>, A<sub>5</sub>, A<sub>6</sub>)

C<sub>3</sub> (1.5, 3.5) (A<sub>2</sub>, A<sub>7</sub>)



# K-medoid

5

$$x_1 \rightarrow (2, 6)$$

$$x_2 \rightarrow (3, 4)$$

$$x_3 \rightarrow (3, 8)$$

$$x_4 \rightarrow (4, 7)$$

$$x_5 \rightarrow (6, 2)$$

$$x_6 \rightarrow (6, 4)$$

$$x_7 \rightarrow (7, 3)$$

$$x_8 \rightarrow (7, 4)$$

$$x_9 \rightarrow (8, 5)$$

$$x_{10} \rightarrow (7, 6)$$

1) Initialize K-centers:

$$\text{let } K = 2$$

$$c_1 = (3, 4), \quad c_2 = (7, 4)$$

<u>i</u>	<u>Cost distance</u>
1	3
2	3
3	4
4	4
5	5
6	3
7	5
8	6
9	6
10	6

$$\text{Cluster 1} = \{ (2, 6), (2, 6), (3, 8), (4, 7) \}$$

for  $C_2$  (7,4)

(6)

$i$	Cost distance
1	7
3	8
4	6
5	3
6	1
7	1
9	2
10	2

Cluster 2  $\rightarrow \{ (7,4), (6,2), (6,4), (7,3), (8,5), (7,6) \}$

So the total cost is  
 $\Rightarrow 2+4+4+3+1+1+2+2$   
 $= 20.$

Step 2: Select one of the nonmedioids ( $O'$ )

let us assume  $O' = (7,3)$

$\therefore$  Now medioids are  $C_1 (3,4)$  &  $O' (7,3)$

$$\text{total cost} = 22 > \underline{20}$$

So moving towards  $O'$  is bad idea.

# Hierarchical Clustering:

(7) (1) (1)

⇒ is used to group similar objects into clusters.

⇒ It is type of unsupervised clustering methods

Step 1 — Given  $n \times n$  distance matrix / Proximity matrix

	A	B	C	D	E	F
A	0	662	887	255	412	996
B	662	0	295	468	268	400
C	887	295	0	754	564	138
D	255	468	754	0	219	869
E	412	268	564	219	0	669
F	996	400	<span style="border: 1px solid black;">138</span>	869	669	0

- 1) Single linkage Method (min)
- 2) Complete - linkage (max)
- 3) Average - linkage (avg.)

It called  
→ Agglomerative  
Hierarchical Clustering  
bcz it merges  
clusters iteratively.

→ Divisive  
↳ rarely used



Use Min/Single Linkage Method.

(8)

step-2

	A	B	C/F	D	E
A	0	662	877	255	412
B	662	0	295	468	268
C/F	877	295	0	754	564
D	255	468	754	0	219
E	412	268	564	219	0

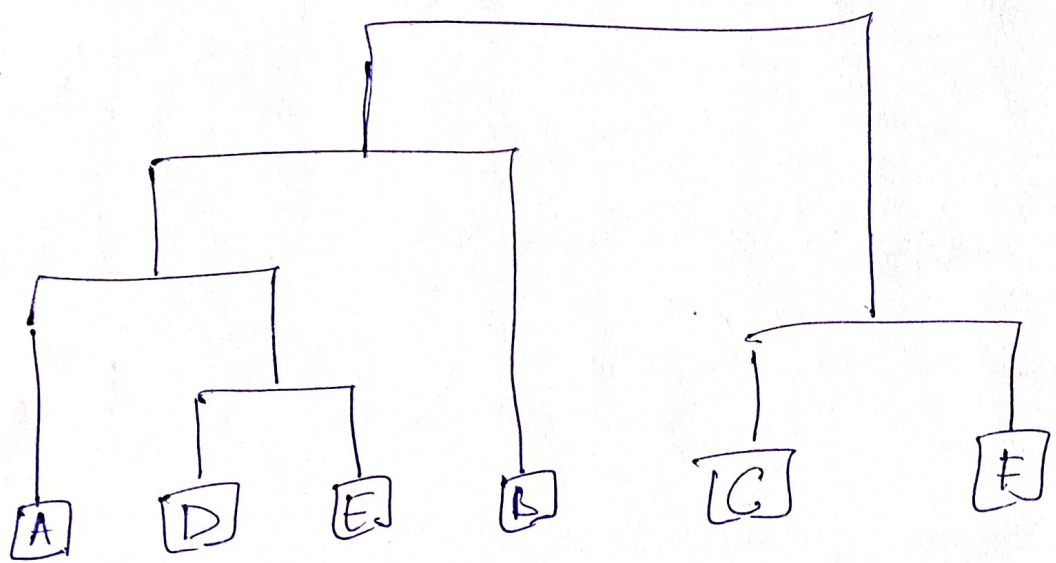
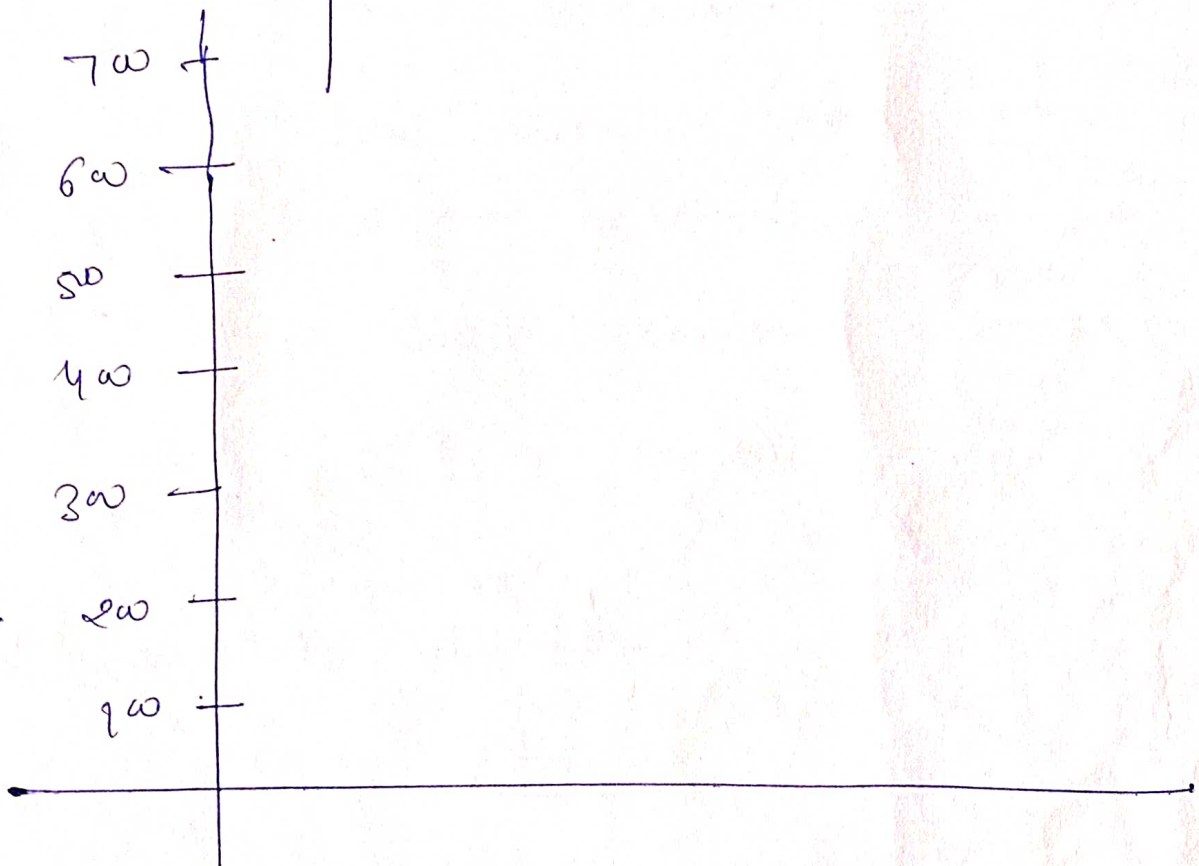
Step-3:

	A	B	C/F	D/E
A	0	662	877	255
B	662	0	295	268
C/F	877	295	0	564
D/E	<del>412</del> 255	<del>268</del> 268	- 564	0

step-4:

	A/D/E	B	C/F
A/D/E	0	<del>662</del> 268	<del>877</del> 564
B	<del>662</del> 268	0	295
C/F	564	295	0

	A/B/D/E	C/F	
A/B/D/E	0	295	(9)
C/F	295	0	

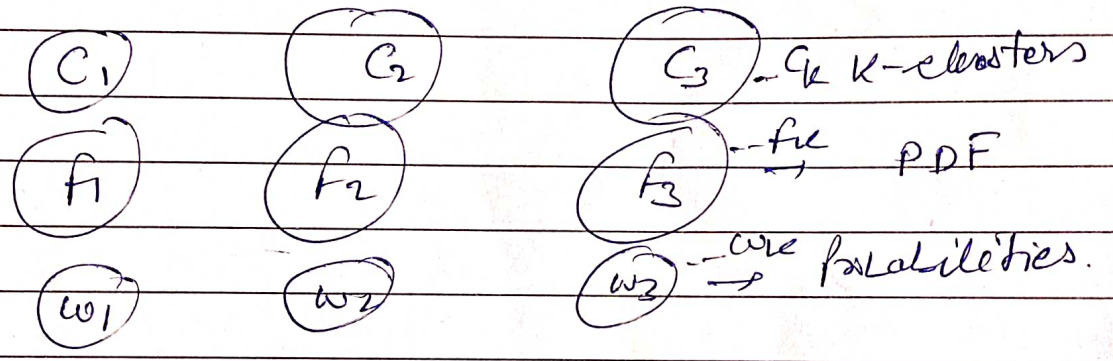


Gaussian Mixture Model:

Given Data set  $\rightarrow D$

$K \rightarrow$  No. of clusters

The task of Probabilistic model based cluster analysis is to infer a set of probabilistic clusters that is most likely to generate D.



$\Rightarrow$  For an object  $O$ , the probability that  $O$  is generated by cluster  $C_j$  ( $1 \leq j \leq k$ ) is given by

$$P(O|C_j) = w_j f_j(O) \rightarrow (1)$$

Therefore the probability that  $O$  is generated by the set  $C$  clusters is

$$P(O|C) = \sum_{j=1}^k w_j f_j(O) \rightarrow (2)$$

DATE: / /

Teacher Signature \_\_\_\_\_



⇒ For a data set  $D = \{o_1, o_2, o_3, \dots, o_n\}$  of  $n$ -objects, we have

$$P(D|C) = \prod_{i=1}^n P(o_i | C) = \prod_{i=1}^n \sum_{j=1}^K w_j P_j(o_i | \theta_j) \quad (3)$$

⇒ Now it is clear that the task of probabilistic model-based cluster analysis on a data set  $D$ , is to find a set  $\mathcal{C}$  of  $K$  probabilistic clusters such that  $P(D|C)$  is maximized.

⇒ Let  $o_1, o_2, \dots, o_n$  be  $n$  observed objects &  $\theta_1, \theta_2, \dots, \theta_K$  be the parameters of  $K$ -distributions, denoted by

$$\Theta = \{ \theta_1, \theta_2, \dots, \theta_K \}$$

$$\Theta = \{ \theta_1, \theta_2, \dots, \theta_K \} \text{ resp.}$$

⇒ Then for any object  $o_i \in \Theta$  ( $1 \leq i \leq n$ )

eq<sup>n</sup> (2) can be rewritten as

$$P(o_i | \Theta) = \sum_{j=1}^K w_j P_j(o_i | \theta_j)$$



where  $P_j(o_i | \theta_j)$  is the probability that  $o_i$  is generated from  $j^{th}$  distribution using parameter  $\theta_j$ ,  $\text{eqn (3)}$  can be rewritten as

$$P(O|\theta) = \prod_{i=1}^n \sum_{j=1}^k w_j P_j(o_i | \theta_j)$$

### Univariate Gaussian Mixture Model:

⇒ Assume that PDF of each cluster follows a 1-D Gaussian distribution.

⇒ Let no. of clusters =  $K$

⇒ Two parameters for the PDF of each cluster are center,  $\mu_j$  & std. dev.  $\sigma_j$

( $1 \leq j \leq K$ ) , we denote parameters as  $\theta_j = (\mu_j, \sigma_j)$  &  $\theta = \{ \theta_1, \theta_2, \dots, \theta_K \}$

⇒ Let the data set be  $O = \{ o_1, o_2, \dots, o_n \}$  where  $o_i$  ( $1 \leq i \leq n$ ) is a real number

⇒ For any point  $o_i \in O$ , we have

$$P(o_i | \theta_j) = \frac{1}{\sigma_j \sqrt{2\pi}} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$



⇒ Assuming that each cluster has same probability that  $w_1 = w_2 = \dots = w_k = 1/k$

$$\Rightarrow P(o_i | \theta) = \frac{1}{k} \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} \cdot e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$

$$\Rightarrow P(O | \theta) = \frac{1}{k} \prod_{i=1}^n \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}} \quad \text{--- (4)}$$

⇒ The task of probabilistic model-based cluster analysis using univariate Gaussian mixture model is to infer  $\theta$  such that eq<sup>n</sup> (4) is minimized.

⇒ using the EM algorithm for mixture models:

Given a set of objects  $O = \{o_1, o_2, \dots, o_n\}$  we want to mine a set of parameters  $\theta = \{\mu_1, \mu_2, \dots, \mu_k\}$  such that  $P(O | \theta)$  is maximized.

where  $\theta_j = (\mu_j, \sigma_j)$

$\theta_j$  univariate Gaussian Distribution



EM Algo:

1) We assign random values to parameters  $\theta$  as the initial values, we then iteratively conduct E-step & M-step until the parameters converge or the change is sufficiently small.

E-Step  $\rightarrow$  1) for each object  $o_i \in O (1 \leq i \leq n)$  we calculate probability that  $o_i$  belongs to each distribution that is

$$P(\theta_j | o_i, \theta) = \frac{P(o_i | \theta_j)}{\sum_{k=1}^K P(o_i | \theta_k)}$$

M-Step  $\rightarrow$  We adjust the parameters  $\theta$  so that the expected likelihood  $P(O | \theta)$  in eq<sup>n</sup> (4) is maximized. This can be achieved by setting

$$\mu_j = \frac{1}{K} \sum_{i=1}^n o_i \frac{P(\theta_j | o_i, \theta)}{\sum_{i=1}^n P(\theta_j | o_i, \theta)}$$

$$\mu_j = \frac{1}{K} \sum_{i=1}^n o_i P(\theta_j | o_i, \theta) / \sum_{i=1}^n P(\theta_j | o_i, \theta)$$

$$\sigma_j =$$

$$\frac{\sum_{i=1}^n P(\theta_j | x_i, \theta) (x_i - \mu_j)^2}{\sum_{i=1}^n P(\theta_j | x_i, \theta)}$$

DATE: / /

Teacher Signature



# Probability Distribution Function (PDF)

$$P.D.F \Rightarrow P(a < x < b) = \int_a^b f(x) dx = \int_a^b \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

$x \rightarrow$  Random variable

$\mu \rightarrow$  mean value

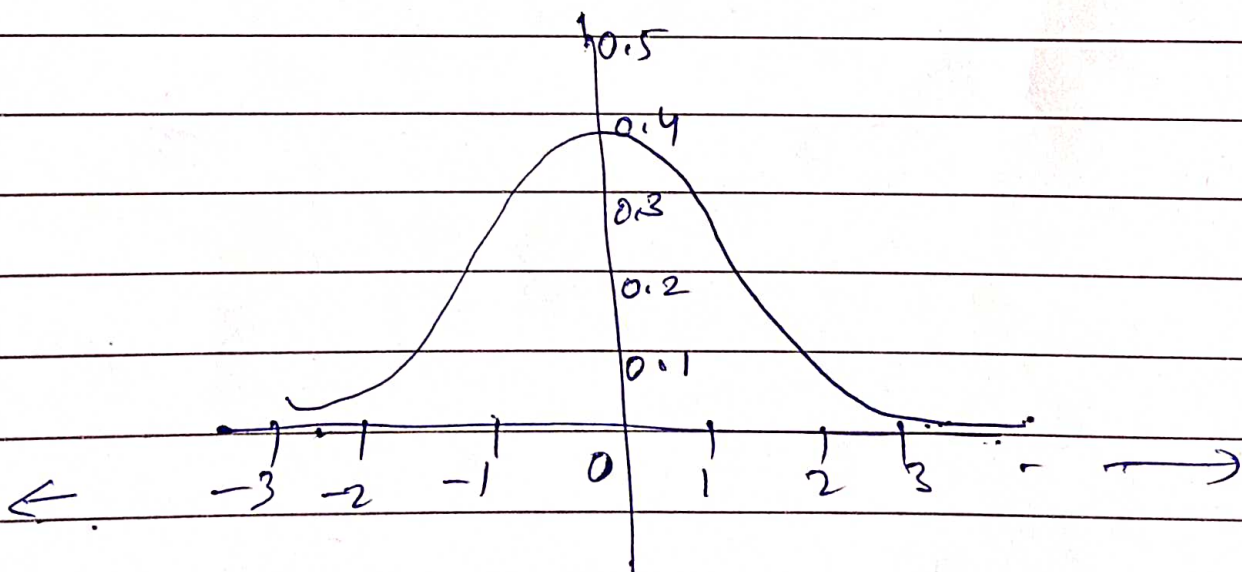
$\sigma \rightarrow$  std. devi.

## Probability Density Function for Normal Distribution

$$f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

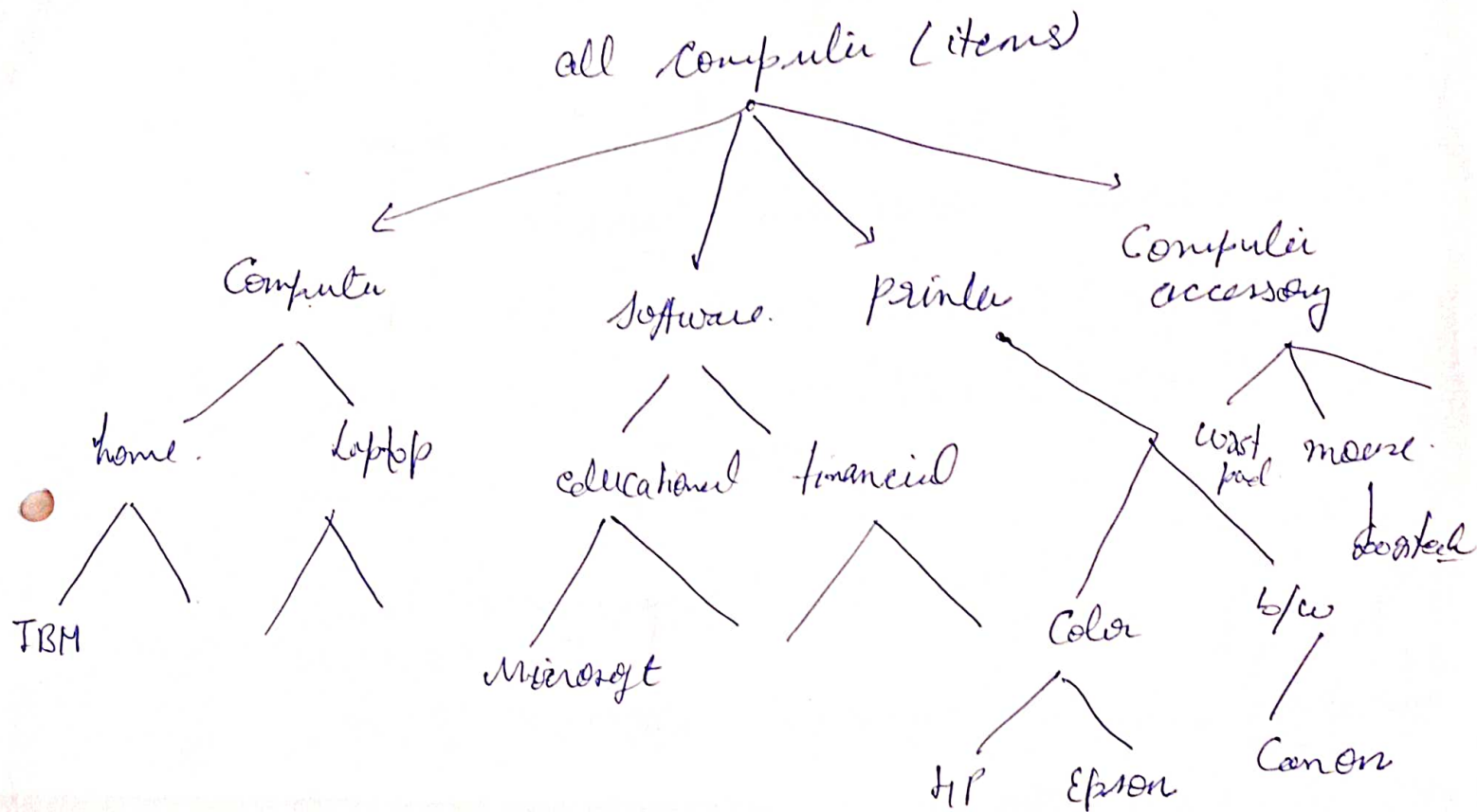
eg  $\mu=0, \sigma=1$

$$f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$





Association Rule:



1)  $\implies$  Single-Dimensional or Intra-Dimensional Association Rules:

- buys (X, "IBM home computer")
- $\implies$  buys (X, "printer")
- $\implies$  (buys (X, "Sony b/w printer")
- $\implies$  buys (X, "b/w printer")
- $\implies$  - - -

2) Multidimensional:

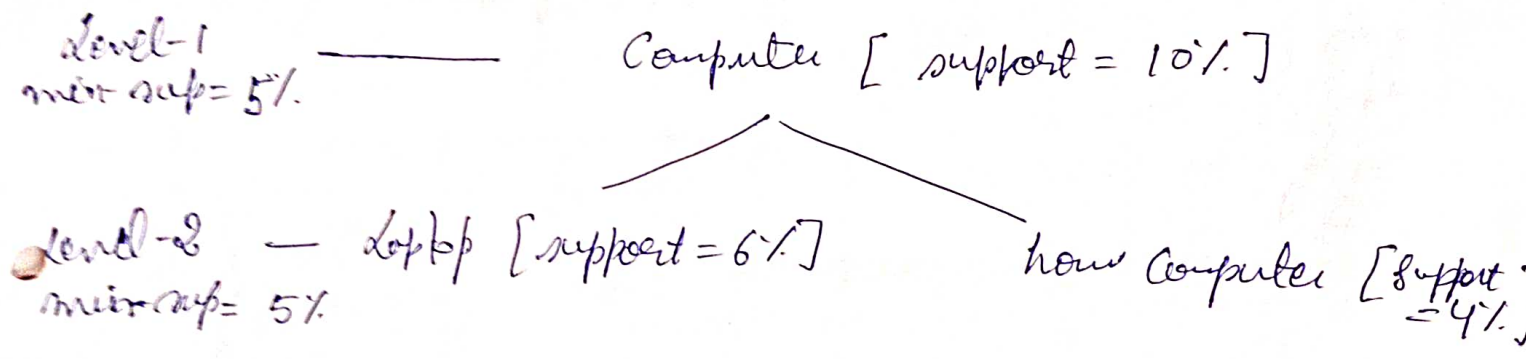
- age (X, "19-24")  $\wedge$  occupation (X, "student")
- $\implies$  buys (X, "laptop")

→ MDAR with no repeated predicates are called: inter-dimension association = rule.

→ with repeated predicates: called ⇒ hybrid-dimension association rules  
age (X, "19-24") ∧ buys (X, "laptop")  
⇒ buys (X, "printer")

Multilevel Mining:

1) Multilevel Mining with uniform support:



2) Multilevel Mining with reduced support:

level-1 [min support = 5%]  
level-2 [min support = 3%]

⇒ Rules generated from mining with concept hierarchies are called multilevel association rules.

⇒ Approaches to Mining multilevel association rules:

1) Using uniform min support for all levels :- (

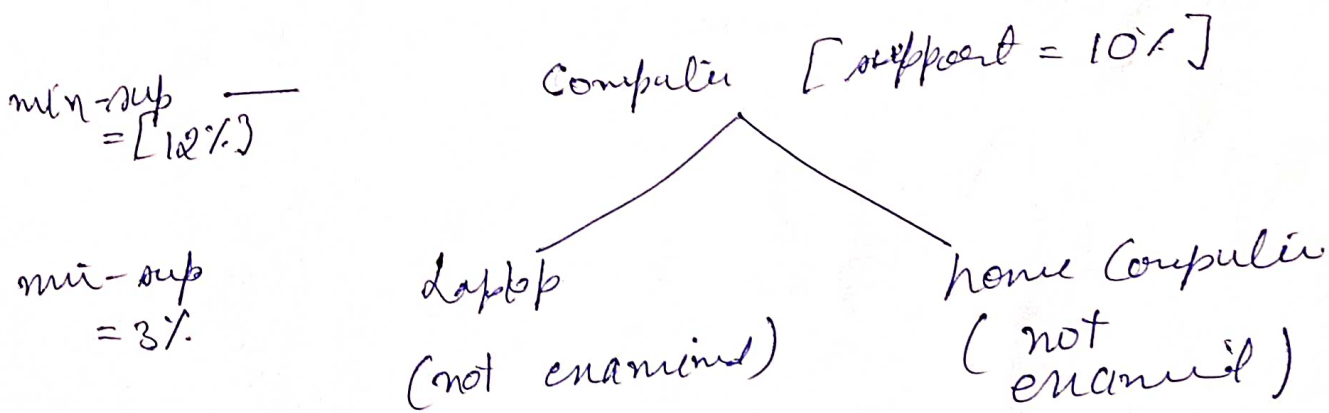
2) Using reduced min support at lower levels:

⇒ For mining MLA with reduced support, there are alternative search strategies are:

a) level-by-level independent:

→ each node is examined, regardless of whether or not its parent node is found to be frequent.

b) level-cross filtering by single item:



⇒ An item at  $i^{th}$  level is examined if & only if its parent node at  $(i-1)$  level is frequent.



c) level-cross filtering by K-itemset. (20)

min-sup = 5%

Computer & printer [support = 7%]

min  
supp  
= 2%

laptop Computer  
& b/w printer  
[support = 1%]

laptop Computer  
& color printer  
[support = 2%]

home Computer  
& b/w printer  
[support = 1%]

home Computer &  
color printer  
[support = 8%]

⇒ A K-itemset at the  $i^{\text{th}}$ -level is examined if & only if its corresponding parent K-itemset at the  $(i-1)^{\text{th}}$  level is frequent.

⇒ Computer & printer is frequent

- ⇒ same concept level.
- ⇒ different concept level.  
(cross-level association rule)

⇒ How can ~~cross-association~~ cross-level - association be mined?  
level - i  
level - j → reduced. min support

threshold of level j should be used overall so that items from level j can be included in the analysis.

⇒ checking for Redundant multilevel association rules?

⇒ concept hierarchies are useful in data mining since they permit the discovery of knowledge at different levels of abstraction, such as multilevel association rules.

- R1 ⇒ home computer ⇒ b/w printer [s = 8%, c = 70%]
- R2 ⇒ IBM home computer ⇒ 4w printer [s = 2%, c = 74%]



# ⇒ Multidimensional association rules :

(22)

- 1) Single D or Intra-D
- 2) HD. - with no repeated predicates (Inter-D)  
- with repeated predicates (Hybrid-D)

Ex. age (X, "19-24") ∧ buys (X, "laptop")  
⇒ buys (X, "b/w printer")

⇒ database attribute can be  
⇒ Categorical & Quantitative

⇒ Techniques for mining MDAR can be categorized according to three basic approaches regarding the treatment of quantitative attributes:

- 1) Quantitative attributes are discretized using predefined concept hierarchies  
→ This discretization occurs prior to mining

Ex: Concept hierarchy for income may be used to replace the original numeric values of this attribute by ranges "0-20K", "21-30K" & so on.

→ Here discretization is static & predetermined  
⇒ Mining MDAR using static discretization  
of quantitative attributes.

2)

## Quantitative AR:

Quantitative attributes are discretized into "bins" based on distribution of data.

⇒ These bins may be further combined during mining process. The discretization process is dynamic.

3)

Distance-Based AR: Quantitative attributes are discretized so as to capture the semantic meaning of neighboring interval data. This dynamic discretization procedure considers the distance b/w data points.



# Mining Frequent Itemsets:

(Apriori Algorithm)

Introduction: - In data mining, association rule learning is popular & well researched method for discovering interesting relations b/w variables in large databases.

Ex. rule { onion, potato }  $\Rightarrow$  { burger }  
{ milk, bread }  $\Rightarrow$  { butter }  
 $\Rightarrow$  In Computer Science & Engg. data mining,

Apriori is a classic algorithm for learning association rules.

Apriori Algo.  $\rightarrow$  is an influential algo. for mining frequent itemsets for boolean association rules.

- 1) Frequent Item sets: - The sets of item which has min support ( denoted by  $L_k$  or  $i_k$  - itemset)
- 2) Apriori Property - Any subset of frequent itemset must be frequent.
- 3) Join operation - To find  $L_k$ , a set of candidate  $k$ - itemsets is generated by joining  $L_{k-1}$  with itself.



steps:- 1) Find the frequent items:  
the sets of items that have min support.

→ A subset of a frequent itemset must also be a frequent itemset.

i.e. if  $\{A, B\}$  is a frequent itemset, both  $\{A\}$  &  $\{B\}$  should be frequent itemset.

- 2) Iteratively find frequent itemsets with Cardinality from 1 to K. (K-itemset)
- 3) Use the frequent itemsets to generate association rule.

Pseudo-Code:-

$C_k$ : Candidate itemset of size k.

$L_k$ : frequent itemset of size k.

$L_1 = \{ \text{frequent items} \}$

for (  $k=1$  ;  $L_k \neq \emptyset$ ,  $k++$  ) do begin

$C_{k+1} =$  candidates generated from  $L_k$ .

for each transaction  $t$  in database do increment the counter of all candidates in  $C_{k+1}$  that are contained in  $t$ .

$L_{k+1} =$  candidates in  $C_{k+1}$  with min-support  
end return  $\cup_k L_k$ ;

Ex:

TID	list of items
T100	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>
T100	I <sub>2</sub> , I <sub>4</sub>
T100	I <sub>2</sub> , I <sub>3</sub>
T100	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>
T100	I <sub>1</sub> , I <sub>3</sub>
T100	I <sub>2</sub> , I <sub>3</sub>
T100	I <sub>1</sub> , I <sub>3</sub>
T100	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>
T100	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>

11/9

⇒ Total 9- transactions in database D.

⇒ Suppose min support count is 2.  
∴  $(2/9) = 22\%$

⇒ Let min confidence required is 70%.

⇒ find out the freq. itemset using apriori algo.

⇒ Then association rules will be generated using min. support & min. confidence.

Step-1: Generating 1-itemset Frequent Pattern

Item	supp count
{I <sub>1</sub> }	6
{I <sub>2</sub> }	7
{I <sub>3</sub> }	6
{I <sub>4</sub> }	2
{I <sub>5</sub> }	2

C<sub>1</sub>

⇒ Compare Candidate support count with min support count.

⇓

Itemset	Sup-Count
{I <sub>1</sub> }	6
{I <sub>2</sub> }	7
{I <sub>3</sub> }	6
{I <sub>4</sub> }	2
{I <sub>5</sub> }	2

L<sub>1</sub>

(set of frequent 1-itemset)



Step 2 : generating 2-itemset freq. patterns

Generate  $C_2$  candidates from  $L_1 \rightarrow$

Itemset	Supp. Count
$\{I_1, I_2\}$	4
$\{I_1, I_3\}$	4
$\{I_1, I_4\}$	1
$\{I_1, I_5\}$	2
$\{I_2, I_3\}$	4
$\{I_2, I_4\}$	2
$\{I_2, I_5\}$	2
$\{I_3, I_4\}$	0
$\{I_3, I_5\}$	1
$\{I_4, I_5\}$	0

$C_2$

Compare candidate support count with min support count

$I_1 I_2$   
 $I_1 I_3$   
 $I_1 I_5$   
 $I_2 I_3$   
 $I_2 I_4$   
 $I_2 I_5$

$L_2$

Itemset	Supp. Count
$\{I_1, I_2\}$	4
$\{I_1, I_3\}$	4
$\{I_1, I_5\}$	2
$\{I_2, I_3\}$	4
$\{I_2, I_4\}$	2
$\{I_2, I_5\}$	2

Step 3: Generating 3-itemset freq. Pattern. (29)

Itemset	Sup-Count	Itemset	Sup-Count	
$\{I_1, I_2, I_3\}$	2	$\Rightarrow$	$\{I_1, I_2, I_3\}$	2
$\{I_1, I_2, I_5\}$	2		$\{I_1, I_2, I_5\}$	2
		<hr/>		
		$L_3$		

$\Rightarrow$  The generation of the set of candidate 3-itemset  $C_3$ , involves use of the Apriori property.

$\rightarrow$  In order to find  $C_3$ , we compute  $L_2$  join  $L_2$ .

$$\Rightarrow C_3 = L_2 \text{ join } L_2 = \left\{ \begin{array}{l} \{I_1, I_2, I_3\}, \{I_1, I_2, I_5\} \\ \{I_1, I_2, I_4\}, \{I_1, I_3, I_5\}, \{I_1, I_3, I_4\} \\ \{I_2, I_3, I_4\}, \{I_2, I_3, I_5\}, \{I_2, I_4, I_5\} \end{array} \right\}$$

$\Rightarrow$  Now, join step is complete and Prune step will be used to reduce the size of  $C$ . Prune step helps to avoid heavy computation due to large  $C_k$ .



⇒ Based on Apriori Algo. Property that all subsets of a freq. itemset must also be frequent. ~~we can determine that.~~

For Example: Let take  $\{I_1, I_2, I_3\}$  ⇒ The 2-item subset of it are  $\{I_1, I_2\}$ ,  $\{I_1, I_3\}$  &  $\{I_2, I_3\}$ .  
Since all 2-item subsets of  $\{I_1, I_2, I_3\}$  are members of  $L_2$ , we will keep  $\{I_1, I_2, I_3\}$  in  $C_3$ .

⇒ Let's take another example of  $\{I_2, I_3, I_5\}$   
subset:  $\{I_2, I_3\}$ ,  $\{I_3, I_5\}$ , &  $\{I_2, I_5\}$   
↓  
is not in  $L_2$ , hence it is not frequent, Thus we will have to remove  $\{I_2, I_3, I_5\}$  from  $C_3$ .

⇒ Therefore:  $C_3 = \{ \{I_1, I_2, I_3\}, \{I_1, I_2, I_5\} \}$



Step 4: Generating 4-itemset frequent pattern.

$\Rightarrow C_4 = L_3 \text{ join } L_3$

$= \{ \{ I_1, I_2, I_3, I_5 \} \}$

$\Rightarrow$  This itemset is pruned, since its subset  $\{ I_2, I_3, I_5 \}$  is not frequent.

$\Rightarrow$  Thus  $C_4 = \emptyset$  & algo. terminates.

$\Rightarrow$  Now, these frequent itemsets will be used to generate strong association rules (where strong association rules satisfy both min support & min confidence)

Step 5: Generating Association Rules from Frequent Itemsets:

$\Rightarrow$  For each frequent itemset "I" generate all nonempty subsets of I.

$\Rightarrow$  For each nonempty subset S of I, output the rule "S  $\rightarrow$  (I-S)" if  $\text{support-count}(I) / \text{support-count}(S) \geq \text{min-conf.}$

Ex:

We had  $L = \{ \{I_1\}, \{I_2\}, \{I_3\}, \{I_4\}, \{I_5\}, \{I_1, I_2\}, \{I_1, I_3\}, \{I_1, I_5\}, \{I_2, I_3\}, \{I_2, I_4\}, \{I_2, I_5\}, \{I_1, I_2, I_3\}, \{I_1, I_2, I_5\} \}$

Let's take  $I = \{I_1, I_2, I_5\}$

⇒ Its all nonempty subsets are:

$\{ \{I_1, I_2\}, \{I_1, I_5\}, \{I_2, I_5\}, \{I_1\}, \{I_2\}, \{I_5\} \}$

⇒ Let min. confidence is 70%

⇒ The resulting association rules are:

⇒  $R_1: I_1 \wedge I_2 \rightarrow I_5$

⇒ Confidence =  $\frac{sc\{I_1, I_2, I_5\}}{sc(I_1, I_2)}$   
 $= \frac{2}{4} = 50\%$

$R_1$  is rejected.

⇒  $R_2: I_1 \wedge I_5 \rightarrow I_2$

Confidence =  $\frac{sc\{I_1, I_2, I_5\}}{sc(I_1, I_5)}$   
 $= \frac{2}{2} = 100\%$

$R_2$  is selected.

$$R_3: I_1 \rightarrow I_2 \wedge I_5$$

$$\text{confidence} = \frac{\text{sc} \{ I_1, I_2, I_5 \}}{\text{sc} \{ I_1 \}} \\ = \frac{2}{6} = 33\%$$

$R_3$  is rejected.

$$R_4: I_2 \rightarrow I_1 \wedge I_5$$

$$\text{confidence} = \frac{\text{sc} \{ I_1, I_2, I_5 \}}{\text{sc} \{ I_2 \}} \\ = \frac{2}{7} = 29\%$$

$R_4$  is rejected

$$R_5: I_5 \rightarrow I_1 \wedge I_2$$

$$\text{Confidence} = \frac{\text{sc} \{ I_1, I_2, I_5 \}}{\text{sc} \{ I_5 \}} \\ = \frac{2}{2} = 100\%$$

$R_5$  is selected;

$\Rightarrow$  In this way we can found strong association rules.



FP Growth

SC=3

(34)

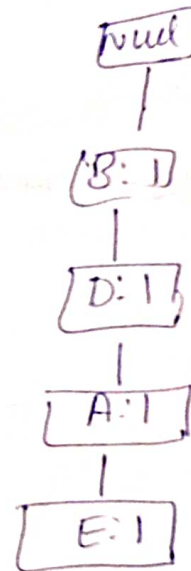
TID	Items
1	E, A, D, B
2	D, A, C, E, B
3	C, A, B, E
4	B, A, D
5	D
6	D, B
7	A, D, E
8	B, C

A → 5 → 3  
 B → 6 → 1  
 C → 3 → 5  
 D → 6 → 2  
 E → 4 → 4

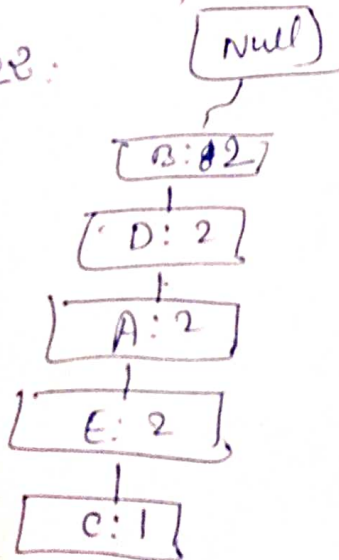
B:6, D:6, A:5, E:4, C:3

TID	Items	Ordered Items
1	E, A, D, B	B, D, A, E
2	D, A, C, E, B	B, D, A, E, C
3	C, A, B, E	B, A, E, C
4	B, A, D	B, D, A
5	D	D
6	D, B	B, D
7	A, D, E	D, A, E
8	B, C	B, C

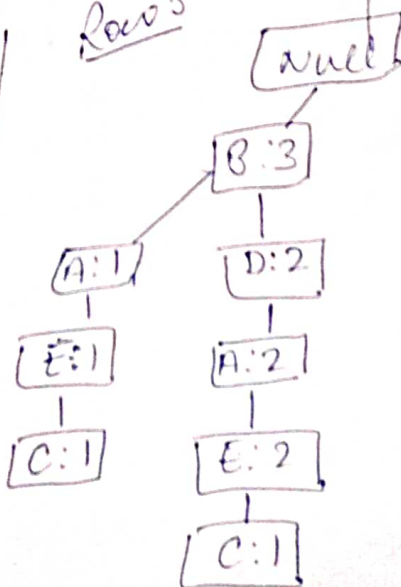
FP-Tree



Row 2:



Row 3



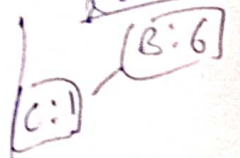
Row 4:

B:4  
D:3  
A:3

Row 5

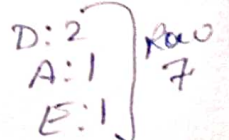


Row 8:



Row 6:

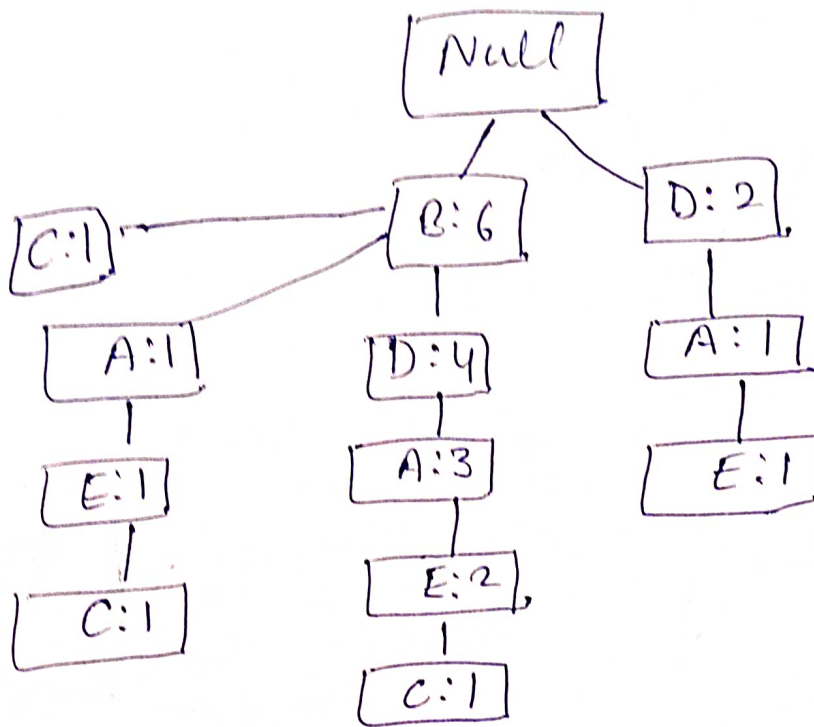
B:5  
D:4



Row 7

# Final Tree

35



# Gaussian Mixture Models:

(26)

1) Probabilistic clustering methods assign instances to clusters probabilistically.

2) The foundation for probabilistic clustering is a statistical model called mixture models.

→ A Mixture model is a set of  $K$ -probability distributions, representing  $K$ -clusters that govern the attribute values for members of that cluster.

→ GMM → Model → is a set of Gaussian (or normal) distributions with mean & covariances.

→ The clustering problem is to take a set of instances & a predefined no. of clusters (Gaussian distributions), & work out each cluster's mean & covariance, & population distribution for the clusters.

→ Let  $x = \{x_j; j = 1, 2, \dots, n\}$  be an  $n$ -D vector to be modelled using Gaussian mixture distribution

→ Let us assume that the model has  $K$ -sub classes (clusters).

→ Then the following parameters are required to completely specify the  $k^{\text{th}}$  subclass:  $k = 1, \dots, K$

→  $w_k$  → the probability that a data sample  $x^{(i)}$  belongs to subclass  $k$   
 $i = 1, 2, \dots, n$ , with  $\sum_{k=1}^K w_k = 1$   
 $w_k \geq 0$



→ This parameter gives population distribution b/w clusters;  $w_R = N_R / N$ ; (2)

where  $N_R =$  No. of samples belonging to  $R^{\text{th}}$  cluster.

$\mu_R =$  the  $n$ -D mean vector for subclass  $k$ .

$\Sigma_R =$  the  $n \times n$  covariance matrix for subclass  $k$ .

1) Mixture Model: is a model comprised of an unspecified combination of multiple probability distribution functions.

2) GMM: - is a mixture model that uses a combination of Gaussian (Normal) probability distribution and requires the estimation of the mean & std. deviation parameters for each.

⇒ There are many techniques for estimating the parameters of GMM, although a maximum likelihood is most common.

3) Expectation Maximization algo: - or EM algo.

is an approach for non. likelihood estimation in the presence of latent variables.

4) Latent variables: - <sup>missing</sup> / hidden / unobserved variables.

E-Step: - Estimate the expected value for each latent variable.

M-Step: - optimize the parameters of the distribution using maximum likelihood.



5) Max. likelihood :- Estimation involves treating the problem as an optimization or search problem, where it seeks a set of parameters that result in the best fit for the joint probability of the data sample.

6) Density Estimation :- involves selecting a probability distribution function & the parameters of that distribution that best explain the joint probability distribution of observed data.

### Example of Gaussian Mixture Model:

Problem :- Dataset whose points are generated from one of two Gaussian Processes. The points are 1-D.

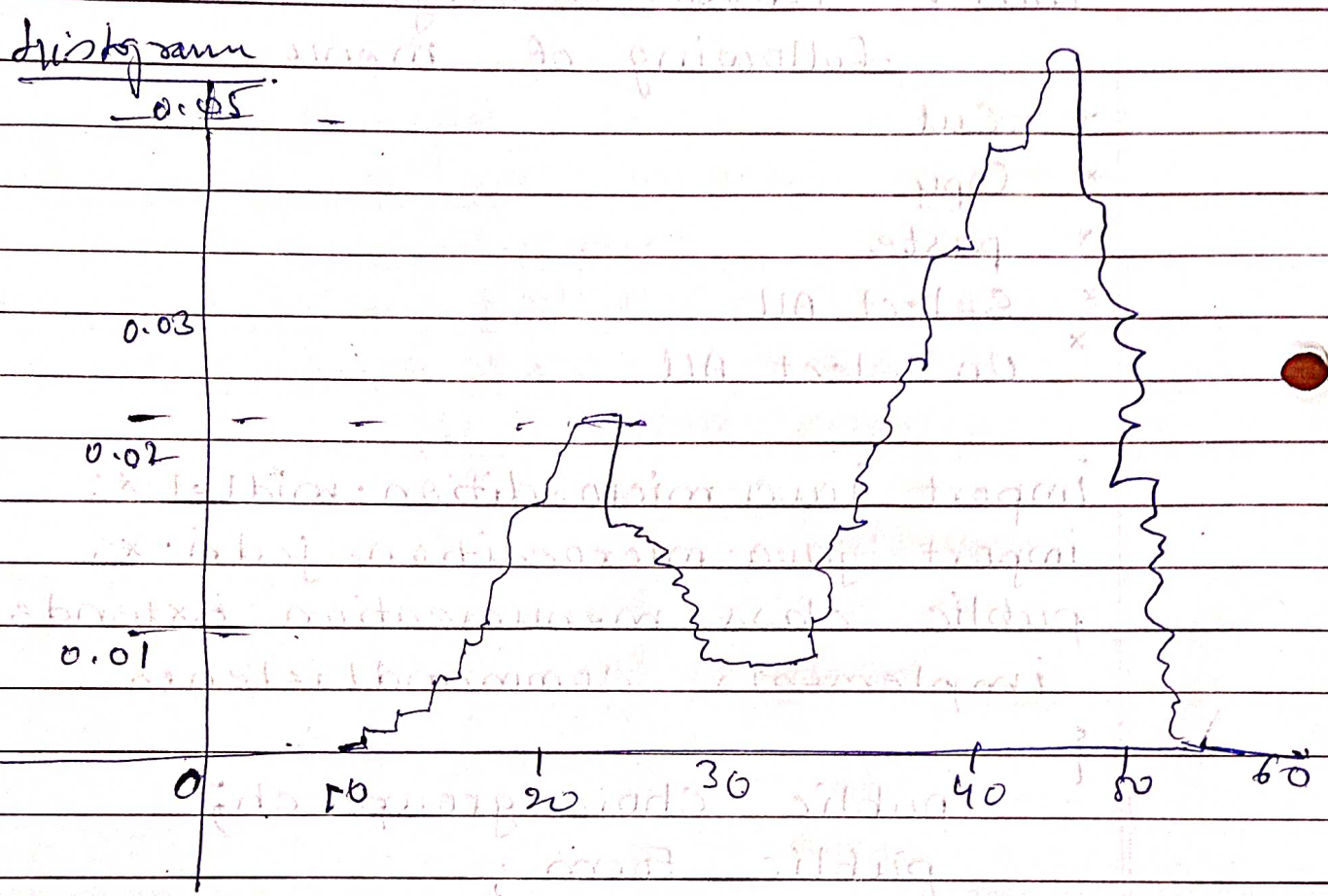
The mean of the first distribution is  $\mu_1 = 20$ .  
second  $\mu_2 = 40$

both distributions have same  
std. dev. = 5

first process  $\rightarrow$  no. of points  $\rightarrow 3000$   
second process  $\rightarrow$  no. of points  $\rightarrow 7000$



⇒ Mix them together.



∴ ⇒ Many of the points in the middle of two peaks that it is ambiguous as to which distribution they were drawn from.

⇒ Model the problem of estimating the density of this dataset via GMM.

⇒ GMM → can be used to model this problem & estimate the parameters of the distributions using EM- Algo.



# EM Coin Flip Example.

H T T T H H T H T H

H H H H T H H H H H

H T H H H H H T H H

H T H T T T H H T T

T H H H T H H H T H

- 1) Initialize  $\theta_A$  &  $\theta_B$  to  
Set  $\theta_A = 0.6$   
 $\theta_B = 0.5$

2) Compute a probability distribution of possible completions of the data using current parameters.

Set 1: H T T T H H T H T H

What is the probability that observe 5 heads & 5 tails in coin A & coin B

⇒ Compute likelihood of set 1 seeing from coin A or B using Binomial distribution with mean probability  $\theta$  on  $n$ -trials with  $k$ -success.

$$P(k) = \binom{n}{k} \theta^k (1-\theta)^{n-k}$$

⇒ likelihood of A → 0.00079

" " B → 0.00097