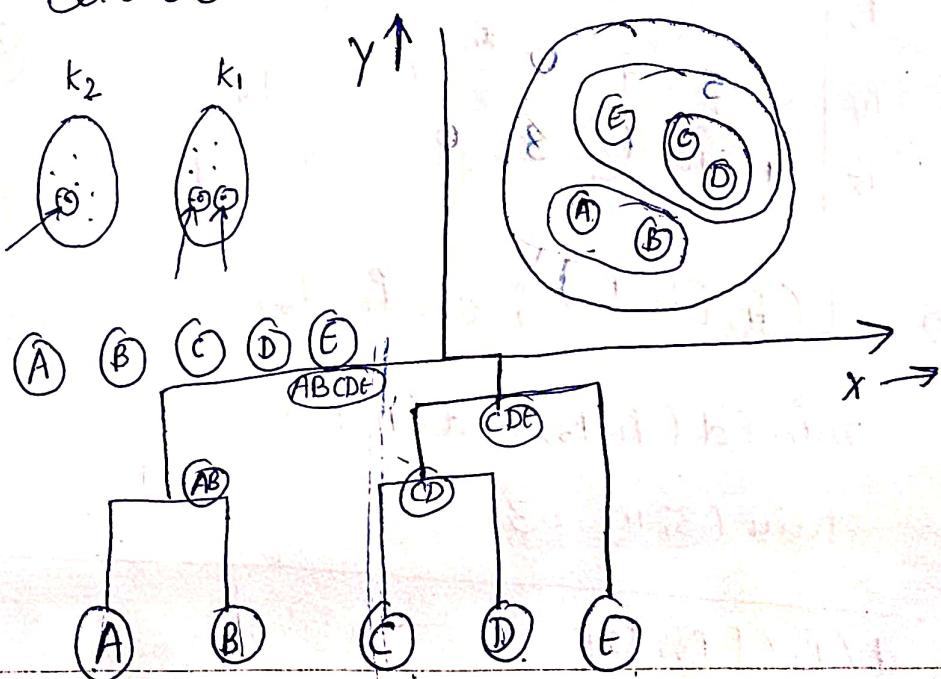


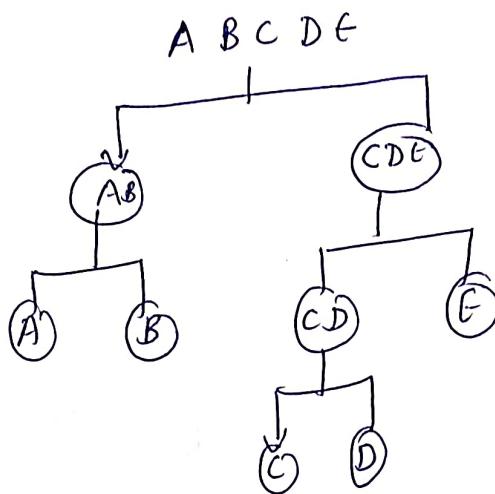
Regression is a statistical that helps us understand and predict the relationship between variables

Hierarchical clustering

- Agglomeration [Dendrogram]
- Divide



Aggns bottom to Top.
 → Divide Algo: - Top to bottom.



Agglomerative clustering :-

	P_1	P_2	P_3	P_4	P_5		P_1	P_2	$\{P_3, P_5\}$	P_4
P_1	6						0			
P_2	9	0					9	0		
P_3	3	7	0				3	7	0	
P_4	6	5	9	0			6	5	8	0
P_5	11	10	2	8	0					

$$\Rightarrow \underline{d(P_1, \{P_3, P_5\})}, \underline{d(P_1, P_5)}$$

$$\min(d(P_1, P_3), d(P_1, P_5))$$

$$\min(3, 11) = 3$$

$$d(P_2, \{P_3, P_5\})$$

$$\Rightarrow \min(d(P_2, P_3), d(P_2, P_5))$$

$$\Rightarrow \min(7, 10) = 7$$

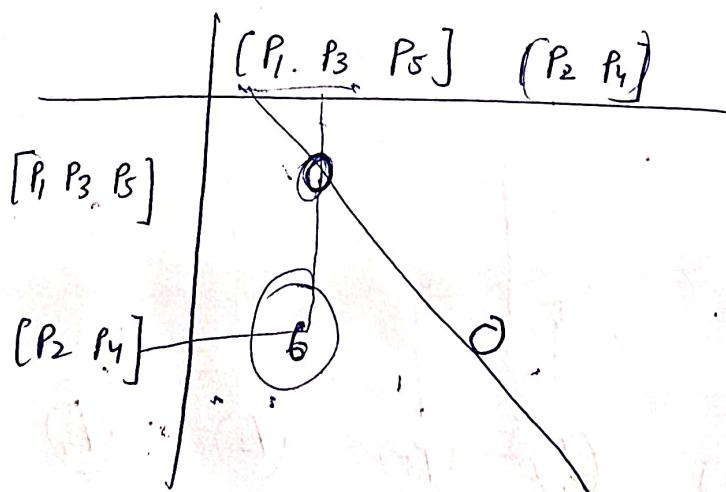
$$d(P_4(P_1 P_3 P_5))$$

$$\Rightarrow \min(P_4, P_1), d(P_4, P_3), d(P_4, P_5)$$

$$\Rightarrow \min(6, 9, 8)$$

$$\Rightarrow \underline{6} \checkmark$$

	P_1	P_2	P_3	P_4	P_5	
P_1	0					$[P_1, P_3, P_5]$
P_2	9	0				$[P_2, P_4]$
P_3	3	7	0			P_2
P_4	6	5	9	0		P_4
P_5	11	10	2	8	0	



$$[P_1, P_2, P_3] [P_4, P_5]$$

9



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

$$\Rightarrow d(P_4, (P_3, P_5))$$

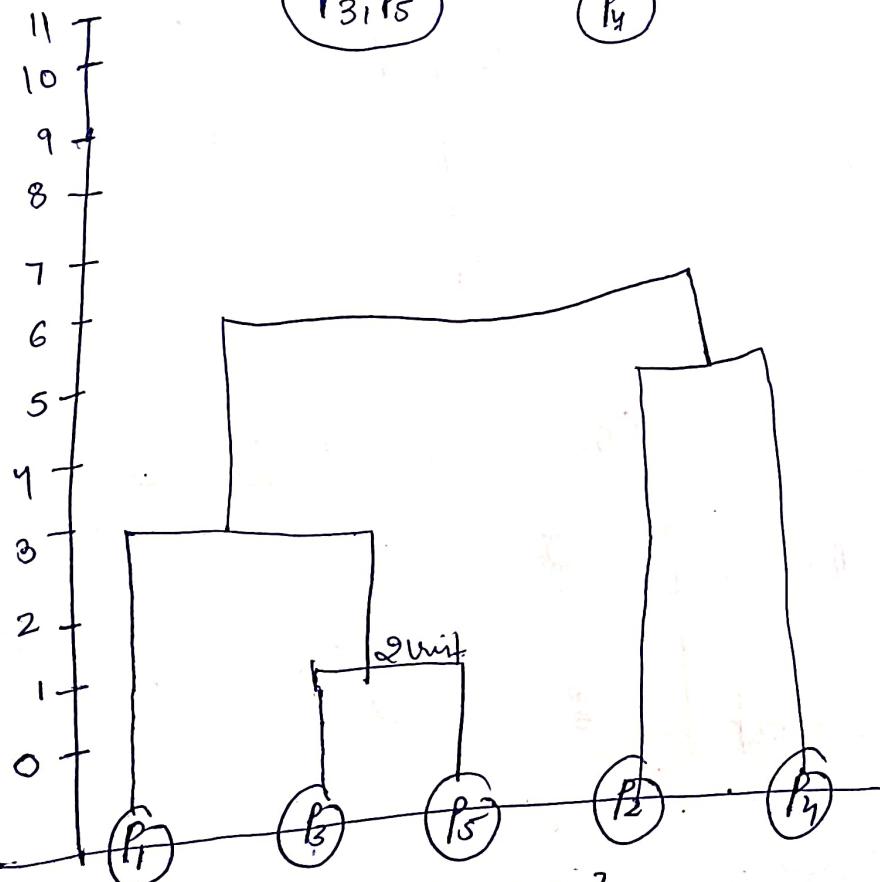
$$\Rightarrow \min d(P_4, P_3), d(P_4, P_5)$$

$$\Rightarrow \min(9, 8) \Rightarrow 8$$



~~[P₁, P₃, P₅] P₂ P₄~~

P₁



$$d(P_2, [P_1, P_3, P_5])$$

$$\min(d(P_2, P_1), d(P_2, P_3), d(P_2, P_5))$$

$$\min(9, 7, 10) = 7$$

~~(P₁, P₃, P₅)~~
(P₂, P₄)

0
6
0

$$\Rightarrow \begin{cases} \min d(P_2, P_1), d(P_2, P_3), d(P_2, P_5), d(P_4, P_1) \\ d(P_4, P_3), d(P_4, P_5) \end{cases}$$

$$\Rightarrow \min(9, 7, 10, 6, 9, 8)$$

$$\Rightarrow 6$$



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

	S_1	S_2	S_3	S_4	S_5	S_6
S_1	0	3.16	2.82	5.09	6.32	(2.23)
S_2		0	3.16	4.47	7.07	5
S_3			0	7.07	8.94	3
S_4				0	(3.16)	7.28
S_5					0	8.06
S_6						0

$\{S_1, S_2, S_3\}$
 $\{S_4, S_5\}$
 $\{S_6\}$
3.17
3.17

Divide Algorithm:- (a, b, c, d, e)

	a	b	c	d	e
a	0	9	3	6	11
b	9	0	7	5	10
c	3	7	0	9	2
d	6	5	9	0	8
e	11	10	2	8	0

$$C_i = \{a, b, c, d, e\}$$



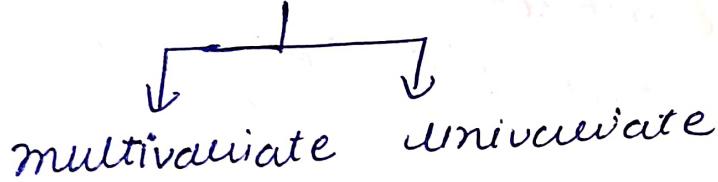
Probabilistic clustering :-

Basic concept of probabilistic model - Based clustering.

mixture models for cluster analysis.

⇒ Gaussian mixture models.

⇒ the expectation - maximization Ealgo.



⇒ analysis of mixture model method.

what is Probabilistic model:-

⇒ model the data from a generative process.

⇒ Assume the data are generated by a mixture of underlying probability distribution.

⇒ Attempt to optimize the fit between the observed data and some mathematical model using a probabilistic approach.

probabilistic model-based clustering :-

- ⇒ each cluster can be represented mathematically by a parametric probability distribution (e.g. a Gaussian or Poisson distribution)
- ⇒ cluster : Data points (or objects) that most likely belongs to the same distribution.
- ⇒ clustering parameter estimation so that they will have a maximum likelihood value to the model by a mixture of k component distribution (e.g. k cluster)

Application:- image segmentation, document clustering.

Two type

mixture model

- ⇒ A very continuous probability distribution is Gaussian distribution it is called as a bell curve. the eq. of one-dimensional

$$P(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

$\sigma \rightarrow$ standard deviation \rightarrow mean

Expectation-maximization algo.

A general technique to find maximum likelihood estimation in mixture models.

EM algo. for Gaussian mixture model.



POORNIMA

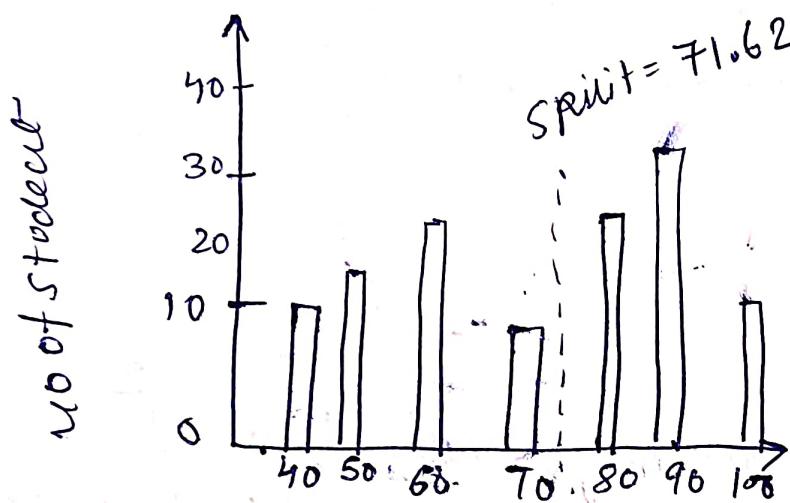
COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

Example:-

PAGE NO.

Histogram of exam score:-



Score	num	Σf_i
40	10	
45	14	
50	19	
55	28	
60	22	
65	19	
70	16	
75	12	
80	30	
85	34	
90	24	
95	19	
100	12	
		$\Sigma f_i = 259$

Step-① split the histogram:-

$$\text{initial split point} = 40 \times 10 + 45 \times 14 + 50 \times 19 \\ + 55 \times 28 + 60 \times 22 + 65 \times 19$$

$$+ 70 \times 16 + 75 \times 12 + 80 \times 30 + 85 \times 34 + \dots$$

$$\text{total mean} : - \frac{90 \times 24 + 95 \times 19 + 100 \times 12}{10 + 14 + 19 + 28 + 22 + 19 + 16 + 12 + 30 + 34 + 24 + 19 + 12} \\ = 71.62$$

Initially, score less than 71.62 is treated as the First component.

Score greater than 71.62 is treated as the Second component.

Step-II initialization:-

First component

$$\underline{\text{mean}} = \frac{40 \times 10 + 45 \times 14 + 50 \times 19 + 55 \times 28 + 60 \times 22 + 65 \times 19 + 70 \times 16}{10 + 14 + 19 + 28 + 22 + 19 + 16}$$

Score	num
40	10
45	14
50	19
55	28
60	22
65	19
70	16

$$m_1 = 56.21$$

$$\underline{\text{variance}} :- 10 \times (40 - 56.21)^2 + 14 \times (45 - 56.21)^2 + \dots + 16 \times (70 - 56.21)^2$$

$$\sigma^2 = \frac{(x_i - \bar{x})^2}{n-1} = \frac{(10+14+19+28+22+19+16)-1}{(10+14+19+28+22+19+16)-1}$$

$$V_1 = 78.64$$

II component:- initialization:-

$$\underline{\text{mean}} = \frac{75 \times 12 + 80 \times 30 + 85 \times 34 + 90 \times 24 + 95 \times 19 + 100 \times 12}{12 + 30 + 34 + 24 + 19 + 12}$$

$$m_2 = 86.68$$

$$\text{variance} = \frac{12 \times (75 - 86.68)^2 + 30 \times (80 - 86.68)^2 + \dots + 12 \times (100 - 86.68)^2}{(12 + 30 + 34 + 24 + 19 + 12) - 1}$$

$$V_2 = 52.16$$



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Step:- 3 Iteration (compute new weight)

this is the eq. to compute the new weight.

$$w_k, \text{new} = \frac{1}{N} \sum_{n=1}^N \frac{w_k \cdot g_1(x_n | m_k, v_k)}{w_k \cdot g_1(x_n | m_k, v_k)}$$

$N \rightarrow$ No. of students

$w_k \rightarrow$ weight (either w_1 or w_2) = 0.5

$m_k - m_1 = 56.21$ & $m_2 = 86.68$

v_k :- variance $v_1 = 78.64$

$v_2 = 52.16$

$x_n \rightarrow$ exam score (40, 45, ..., 100)
from each of the student
259.

$k \rightarrow$ cluster value = 2

$$g_1(x | m_1, v_1) = \frac{1}{(2\pi)^{1/2} |v_1|^{1/2}} \exp \left[-\frac{1}{2} [x - m_1]^T v_1^{-1} [x - m_1] \right]$$

$$= \frac{1}{\text{letsolve numerator First}}$$

$$= \frac{1}{(2 \times 3.14 \times 78.64)^{1/2}} \exp \left[- \frac{(x - 56.21) \times (x - 56.21)}{2 \times (78.64)} \right]$$

$x \rightarrow$ substitute: - 40, 45, 50, 55, 60, 65, 70, 75
80, 85, 90, 95, 100

$$G_1(x_1 | V_1) = -0.008, 0.002, 0.035, 0.045, 0.041, 0.028, \\ 0.013, 0.005, 0.001, 0, 0, 0, 0$$

as II component:-

$$\frac{1}{(2 \times 3.14 \times 52.16)^{1/2}} \exp \left[- \frac{(x - 86.68) \times (x - 86.68)}{2 \times (52.16)} \right]$$

$$x \rightarrow 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100 \\ \downarrow \\ -0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0.001, 0.004, 0.015, 0.036, \\ 0.054, 0.05, 0.028, 0.01$$

$$= \left[\frac{0.5 \times 0.008}{0.5 \times (0.008 + 0)} \times 10 + \frac{0.5 \times 0.002}{0.5 \times (0.002 + 0)} \times 14 + \frac{0.5 \times 0.035}{0.5 \times (0.035 + 0)} \times 19 \right. \\ \left. + \frac{0.5 \times 0.045}{0.5 \times (0.045 + 0)} \times 28 + \frac{0.5 \times 0.041}{0.5 \times (0.041 + 0)} \times 22 \right. \\ \left. + \frac{0.5 \times 0.028}{0.5 \times (0.028 + 0.001)} \times 19 + \dots \right] \times \frac{1}{259} \\ = 0.494$$



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

4

II Substitution:-

$$\left[\frac{0.5 \times 0}{0.5 \times (0.008+0)} \right] \times 10 + \left[\frac{0.5 \times 0}{0.5 \times (0.002+0)} \right] \times 14 + \left[\frac{0.5 \times 0}{0.5 \times (0.035+0)} \right] \times 19 \\ + \left[\frac{0.5 \times 0}{0.5 \times (0.045+0)} \right] \times 28 + \left[\frac{0.5 \times 0}{0.5 \times (0.041 \times 0)} \right] \times 22 \\ - \left[\frac{0.5 \times 0.0}{0.5 \times (0+0.01)} \right] \times 12 \right] \times \frac{1}{259} \\ = 0.506$$

$$w_1 = 0.5$$

$$w_2 = 0.506$$

$$\text{Hence } w_1 = 0.494$$

$$w_2 = 0.506$$

$$w_1 = 0.494 \quad M_1 = 56.45 \quad V_1 = 85.95 \\ w_2 = 0.506 \quad M_2 = 86.46 \quad V_2 = 57.64$$

Em algorithm:- Em algo is a problem in which the data d is a set of instances generated by probability distribution that is mixture of k distinct normal distribution. the case when k=2 and where the distance instances are the point x axis. Each instance is generated using a two-step process.

First one, k normal distribution is selected at random. Second a single random instance x_i is generated according to this selected attributed.

where each of k-normal distribution has a same variance σ^2

the learning task is output μ

$$= \{\mu_1, \dots, \mu_k\}$$

the mean of each of the k distribution.

it is easily calculated max. likelihood hypothesis for the mean of a single normal distribution observed data instance x_1, x_2, \dots, x_n \rightarrow training instance

$$\mu_{ML} = \operatorname{argmin}_u \sum_{i=1}^n (x_i - u)^2$$

in case the sum of squared error is minimized by sample mean.



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

Campus: Course:

Class/Section:
Name of Subject:

Date:
Code:

$$u_{ml} = \frac{1}{M} \sum_{i=1}^m x_i$$

Step I :- calculate the expected value $E[z_i]$ of each hidden variable z_{ij} , assuming the current hypothesis $h = (u_1, u_2)$.

Step 2 :- calculate a new maximum likelihood $h = (u_1, u_2)$ assuming the value taken on by each hidden z_{ij} .

$$E[z_{ij}] = \frac{P(x=x_i) | h=u_j)}{\sum_{n=1}^2 P(x=x_i) | h=u_n)}$$

Association rule mining :-

ARM also called as market

Association rule mining is one of the ways to find pattern in data.

- feature (dimension) which occur together
- feature (ee.) which are correlated.

Association rule in any dataset where feature take only two values of %.

- market Basket analysis is a popular application of association rule mining.

Set of item in transaction is called market Basket.

Example:- For the following given transaction

Data-set generate rule using
Support = 50% confidence = 75%.

measures of effectiveness of the Rule:-

- Support
- Confidence
- lift
- others Affinity, leverage



LURNIMA COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

Campus: Course:

Name of Faculty:

Class/Section:

Name of Subject:

Date:

Code:

Transaction ID	item purchased
1	Bread, cheese, Egg, Juice
2	Bread, milk, cheese, Juice
3	Bread, milk, Yogurt
4	Bread, juice, milk
5	cheese, Juice, milk

SUPPORT

= $\frac{\text{No. of times item appears}}{\text{total no. of transaction}}$

Frequent item set:-

items	Frequency	SUPPORT
✓ Bread	→ 4	$\frac{4}{5} = 80\%$
✓ Cheese	→ 3	$\frac{3}{5} = 60\%$
Egg	→ 1	$\frac{1}{5} = 20\%$
✓ Juice	→ 4	$\frac{4}{5} = 80\%$
✓ Milk	→ 3	$\frac{3}{5} = 60\%$
Yogurt	→ 1	$\frac{1}{5} = 20\%$

minimum SUPPORT 50%

we will remove these two item support is less than 50%.

make 2-item candidate set and wrote their frequency.

item pair	Frequency	SUPPORT
(Bread, cheese) → 2	→ 2/5 = 40%.	
(Bread, juice) → 3	→ 3/5 = 60%.	
(Bread, milk) → 2	→ 2/5 = 40%.	
(cheese, juice) → 3	→ 3/5 = 60%.	
X (cheese, milk) → 1	→ 1/5 = 20%.	
(juice, milk) → 2	→ 2/5 = 40%.	

For Rules → (Bread, juice) — ①

make → (cheese, Juice) — ②

① (Bread, juice)

(Bread → juice) (Juice → Bread)

Confidence:-

$$(A \rightarrow B) = \frac{\text{support of } A \cup B}{\text{support}(A)}$$

① Bread → juice

$$\frac{S(B \cup J)}{S(B)} = \frac{3 \times 5}{5 \cdot 4} = 75\%$$

② Juice → Bread

$$\frac{S(J \cup B)}{S(J)} = \frac{3 \cdot 5}{5 \cdot 4} = 75\%$$

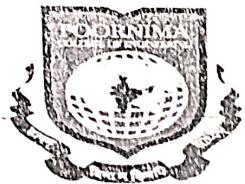
Both rules are good enough to be implemented.

② cheese → juice
Juice → cheese

$$\frac{S(C \cup J)}{S(C)}$$

$$\frac{S(J \cup C)}{S(J)} = \frac{3 \cdot 5}{5 \cdot 3} = 100\%$$

$\frac{3 \cdot 5}{5 \cdot 4} = 75\%$
All rules were covered.



Date → 4/02/2025

PAGE NO.

k-Mean clustering :-

Sr. No	age	amount
c1	20	500
c2	40	1000
c3	30	800
c4	18	300
c5	28	1200
c6	35	1400
c7	45	1800

See Euclidean Distance:-

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

k_1
C1
20, 500

k_2
C2
40, 1000

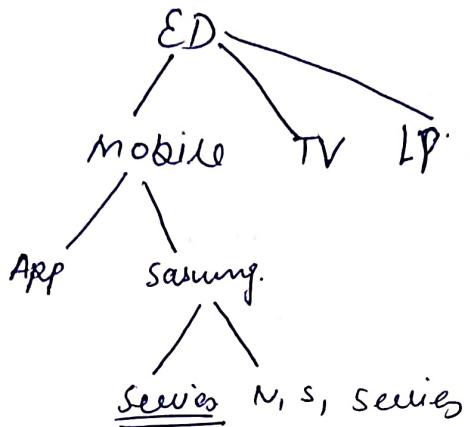
$$C_3 = \sqrt{(30-20)^2 + (800-500)^2} \\ \sqrt{(10)^2 + (300)^2} = \underline{\underline{3000}}$$

$$C_3 = \sqrt{(30-40)^2 + (800-1000)^2} = \underline{\underline{200}}$$

K-mean clustering:-

Central, Hierarchy.

Tree divided.



Step 1 Take mean value.

Step 2 Find nearest no. of mean and put in cluster.

Step 3 :- Repeat one and two until we get same mean.

$$k = \{2, 3, 4, 10, 11, 12, 20, 25, 30\}$$

$k = 2$,

$M_1 = 4$ $M_2 = 12$

k_1
 k_2

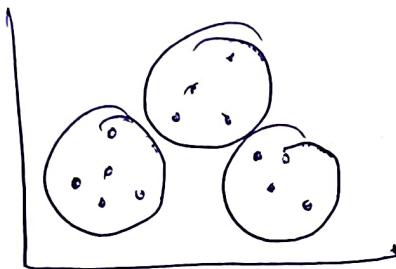
$k_1 = \{2, 3, 4\} \Rightarrow M_1 = 9/3 = 3$

$k_2 = \{10, 11, 12, 20, 25, 30\}$

$M_2 = \frac{108}{6} = 18$

Numerical:-

- k-means clustering is an unsupervised iterative clustering technique.
- it partitions the given data set into k predefined distinct clusters.
- A cluster is defined as a collection of data points certain similarities



$O(nkt)$

$n \rightarrow$ no. of instances

$k \rightarrow$ no. of clusters.

$t \rightarrow$ no. of iterations

$A_1(2, 10), A_2(2, 5), A_3(8, 3), A_4(5, 8)$

$A_5(7, 5), A_6(6, 4), A_7(1, 2), A_8(4, 9)$

initial cluster $A(2, 10), A_4(5, 8)$ and $A_7(1, 2)$

$$P(a, b) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Iteration 01:- $A_1(2, 10), C_1(2, 10)$

$$\begin{aligned} & (2-2) + (10-10) \\ & = 0 \end{aligned}$$

$A_1(2, 10)$ and $C_2(5, 8)$

$$(5-2) + (8-10) = 3+2=5$$



L-means
11
16
Iterat

POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

$$M_1 = 3,$$

$$M_2 = 18$$

PAGE NO.

$$k_1 = \{ 2, 3, 4, 10 \}$$

$$k_2 = \{ 11, 12, 20, 25, 30 \}$$

$$M_1 = 4.75(5), \quad M_2 = 19.6 = (20)$$

$$k_1 = \{ 2, 3, 4, 10, 11, 12 \}$$

$$k_2 = \{ 20, 25, 30 \}$$

$$m_1 = 7, \quad m_2 = 25$$

$$k_1 = \{ 2, 3, 4, 10, 11, 12 \}$$

$$k_2 = \{ 20, 25, 30 \}$$

$$m_1 = 7, \quad m_2 = 25$$

Thus we are getting same measure
have to stop.

New cluster will be

$$K_1 = \{ 2, 3, 4, 10, 11, 12 \}$$

$$K_2 = \{ 20, 25, 30 \}$$

②
fixed
var
value



POORNIMA

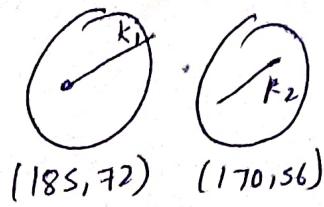
COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

Height weight

PAGE NO.

185	72
170	56
168	60
179	68
182	72
188	77
180	71
186	70
183	84
180	88
180	67
177	76

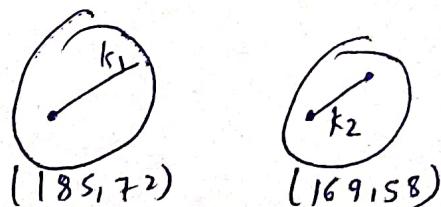


$$k_1 \rightarrow \sqrt{(168-185)^2 + (60-72)^2} \\ = 26.82$$

$$k_2 \rightarrow \sqrt{(168-170)^2 + (60-56)^2} \\ = 4.48$$

New centroid:-

$$k_2 = \left(\frac{170+168}{2} \right), \left(\frac{60+56}{2} \right) = (169, 58)$$



$$k_1 \rightarrow \{1, 4, 5, 6, 7, 9, 10, 11, 12\}$$

$$k_2 \rightarrow \{2, 3\}$$

$$k_1 = \sqrt{(179-185)^2 + (68-72)^2} \\ \Rightarrow \sqrt{(6.82)}$$

$$k_2 = \sqrt{(179-169)^2 + (68-58)^2} \\ = 14.14$$

$A_1(2,10)$, and $(3(1,2) =$

$$|1-2| + |2-10|$$

$$1+8 = 9$$

Given points	Distance $(2,10)$ $(1,2)$	Distance $(5,8)$ $(1,2)$	Distance $(1,2)$ $(4,9)$	Point belongs to class
$A_1(2,10)$	0	5	9	C1
$A_2(2,5)$	5	6	4	C3
$A_3(8,4)$	12	7	9	C2
$A_4(5,8)$	5	0	10	C2
$A_5(7,5)$	10	5	9	C2
$A_6(6,4)$	10	5	7	C2
$A_7(1,2)$	9	10	0	C3
$A_8(4,9)$	3	2	10	C2



DETAILED LECTURE NOTES

Transational ID	items purchased
1	Bread, cheese, egg, juice
2	Bread, cheese, juice
3	Bread, milk, Yogurt
4	Bread, juice, milk
5	cheese, juice, milk

PAGE NO.....

Suppose = 50%.

confidence = 75%.

7
ing

Frequent itemset support \rightarrow Bread

items	Frequency	support
Bread	4	$4/5 \times 100 = 80\%$
cheese	3	$3/5 \times 100 = 60\%$
egg	1	$1/5 = 20\%$
juice	4	$4/5 = 80\%$
milk	3	$3/5 = 60\%$
yogurt	1	$1/5 = 20\%$

item pairs	Frequency	support
(Bread, cheese)	→ 2	$2/5 = 40\%$
(Bread, juice)	→ 3	$3/5 = 60\%$
(Bread, milk)	→ 2	$2/5 = 40\%$
(cheese, juice)	→ 3	$3/5 = 60\%$
(cheese, milk)	→ 1	$1/5 = 20\%$
(juice, milk)	→ 2	$2/5 = 40\%$

For Rules → (Bread, juice)
 → (cheese, juice)

(i) (Bread, juice)

(Bread → juice) (juice → Bread)

$$\text{confidence } (A \rightarrow B) = \frac{\text{support}(A \cup B)}{\text{support } A}$$

(1) Bread → juice = $\frac{s(A \cup B)}{s(A)} = \frac{3/5 \cdot 3/4}{3/5} = 3/4 = 75\%$

(2) juice → Bread = $\frac{3/5}{3/4} = 75\%$

③ $\frac{\text{cheese} \rightarrow \text{juice}}{3/5 \cdot 5/3 = 100\%}$ | $\frac{\text{juice} \rightarrow \text{cheese}}{3/5 \cdot 5/4 = 75\%}$

A. Frequent pattern set (L) is built which will contain all element whose frequency is greater than or equal to the

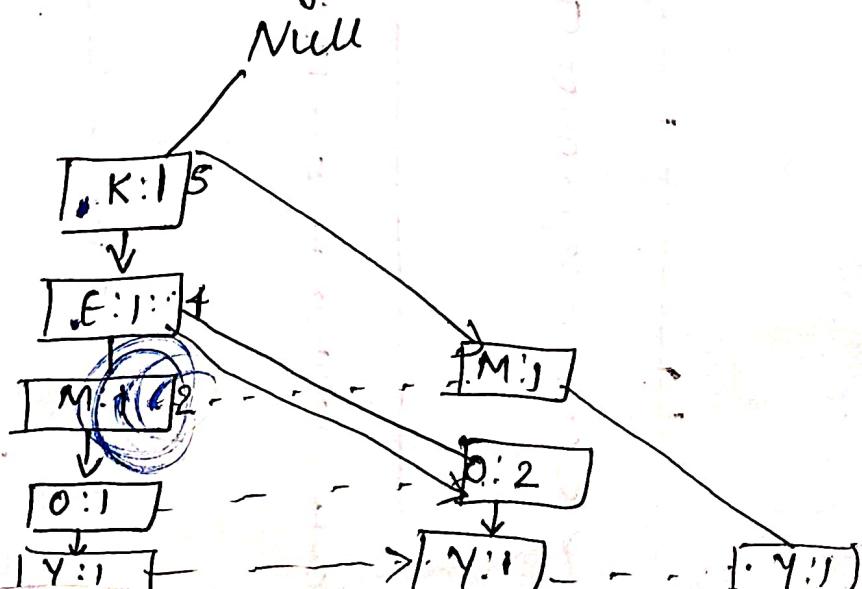
These elements are stored in descending order of their respective.

$$L = \{ K: 5, E: 4, M: 3, O: 3, Y: 3 \}$$

Frequent pattern set $L = \{ K: 5, E: 4, M: 3, O: 3, Y: 3 \}$

Transaction ID	Items	ordered-itemset
T1	{E, K, M, N, O, Y}	{K, E, M, O, Y} {K, E, M, O, Y}
T2	{D, E, K, N, O, Y}	{K, E, O, Y} {K, E, O, Y}
T3	{A, E, K, M}	{A, E, K, M} {K, E, M}
T4	{C, K, M, O, Y}	{C, K, M, O, Y} {K, M, Y}
T5	{C, E, I, K, O, O}	{C, E, I, K, O, O} {K, E, O}

a) inserting the set $\{ K, E, O, Y \}$





DETAILED LECTURE NOTES

PAGE NO.....

Frequent pattern (FP growth):-

Transaction ID	items
T1	{ E, K, M, N, O, Y }
T2	{ D, E, K, N, O, Y }
T3	{ A, E, K, M }
T4	{ C, K, M, O, Y }
T5	{ C, E, I, K, O, O }

the above-given data is a hypothetical dataset of transaction with each letter representing an item

- let the minimum support is 3.

item	Frequency
A	1
C	1
D	2
E	1
I	4
K	1
M	5
N	3
O	2
U	3
Y	3



An autonomous institution approved by RTU, AICTE & UGC • NAAC A+ Accredited

POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

Y	$\{K, E, M, O : 1\}, \{K, E, O : 1\} [K, M : 1]$
O	$\{K, E, M : 1\}, \{K, E : 2\}$
M	$\{K, E : 2\}, \{K : 1\}$
E	$\{K : 4\}$
K	

PAGE NO.....

Y	$\{K, E, M, O : 1\}, \{K, E, O : 1\} [K, M : 1]$	Condition for $\{K : 3\}$
O	$\{K, E, M : 1\}, \{K, E : 2\}$	$\{K, E : 3\}$
M	$\{K, E : 2\}, \{K : 1\}$	$\{K : 3\}$
E	$\{K : 4\}$	$\{K : 4\}$
K		

Y	$\{K, Y : 3\}$
O	$\{K, O : 3\}, \{E, O : 3\}, \{E, K, O : 3\}$
M	$\{K, M : 3\}$
E	$\{E, K : 3\}$
K	



DETAILED LECTURE NOTES

PAGE NO.....

Association rule mining :-

min support

Transactions	items
T1	HOTDOGS, BUNS, KETCHUP
T2	HOTDOGS, BUNS
T3	HOTDOGS, COKE, CHIPS
T4	CHIPS, COKE
T5	CHIPS, KETCHUP
T6	HOTDOGS, COKE, CHIPS

2. Transaction Turnover = 33.34 %.

Q. 2

T - ID	itemset
T = 1000	M, O, N, K, E, Y
T = 1001	D, O, N, K, E, Y
T = 1002	M, A, K, E
T = 1003	M, U, C, K, Y
T = 1004	C, O, O, K, E

Min Support 60%
Condition = 80%.
= Support =
 $\frac{60}{100 \times 5} = 3$

itemset	Supp. count
M	3
O	4
K	5
E	4
Y	3

MO	1
MK	3
ME	2
MY	2
OK	3
O,E	3
O,Y	2
K,E	4
K,Y	3
E,Y	2

→

Now again C₂ in min support

M, K	3
O, K	3
O, E	3
K, Y	3
K, E	4

↓

OKE = 3	MOK 1	1) {Bread, Butter, milk}
OKK = 6	MKE 2	2) {Bread, Butter}
OKE = 5	MKY 2	3) {Bread, cookies, Diaper}
OKE = 5	OKE 3	4) {milk, Diaper, Bread, Butter}
OKE = 5	OKY 2	5) {Bread, Diaper}
OKE = 5	KEO 3	
OKE = 5	KOE 3	
OKE = 5	OKO 3	
OKE = 5	EOK 3	
OKE = 5	OEK 3	
OKE = 5	OKO 3	
OKE = 5	EOO 3	
OKE = 5	OEO 3	
OKE = 5	EOO 3	