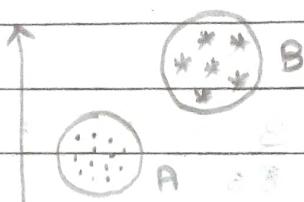


* Clustering

- 1) Partition Based
- 2) Hierarchy Based
- 3) Density Based



Distance Function

$$D(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

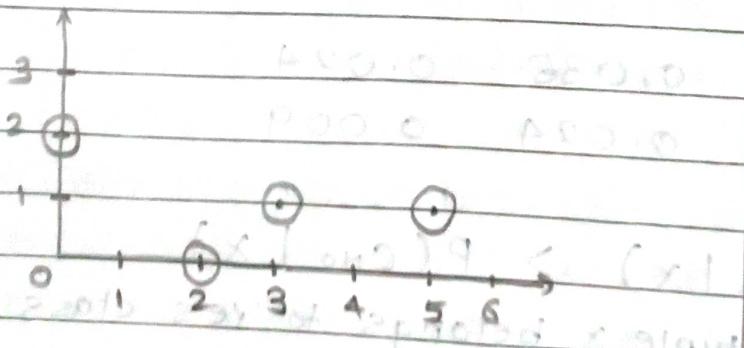
→ Minkowski distance

$$D(x, y) = \left(\sum_{i=1}^n |x_i - y_i| \right)^{1/1}$$

→ Manhattan

$$D(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^2 \right)^{1/2}$$

→ Euclidean



Point	x	y
A	0	2
B	2	0
C	3	1
D	5	1

$$D(P_i, P_j) = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

$$D(A, B) = \sqrt{(0-2)^2 + (2-0)^2} = 2\sqrt{2} = 2.82$$

$$D(A, C) = \sqrt{(0-3)^2 + (2-1)^2} = \sqrt{10} = 3.162$$

$$D(A, D) = \sqrt{(0-5)^2 + (2-1)^2} = \sqrt{26} = 5.099$$

$$D(B, C) = \sqrt{(2-3)^2 + (0-1)^2} = \sqrt{2} = 1.414$$

$$D(B, D) = \sqrt{(2-5)^2 + (0-1)^2} = \sqrt{10} = 3.162$$

$$D(C, D) = \sqrt{(3-5)^2 + (1-1)^2} = \sqrt{94} = 2.82$$

for a given data point convert into Data Matrix

	A	B	C	D
A	0	2.82	3.16	5.09
B	2.82	0	1.41	3.16
C	3.16	1.41	0	2
D	5.09	3.16	2	0

* K-Means Clustering Algorithm (Partition Based)

For given value of K, K-Means is implemented in

4 steps :

- 1) Partition objects into K non-empty subsets.
- 2) Compute the seed points as centroids of the clusters of the current partition. (centroid is the mean point of the cluster)
- 3) Assign each object to cluster with the nearest seed point.

4) Go back to step ②, stop when no more new assignments are there.

Q. For given dataset create 3 clusters using k-means clustering algorithm. (a) (a)

Data = 22, 9, 12, 15, 10, 27, 35, 18, 36, 11

K=3 (a) (a)

→

1st Iteration: (a) (a)

K=1 {22} = {22, 27, 35, 36, 18} Mean = 27.6

K=2 {9} = {9, 10} Mean = 9.5

K=3 {12} = {12, 11, 15} Mean = 12.6

Calculate the mean of each set as the new centroid.

0 9 8 11
PC-2 30-2 18-2 0 11

2nd Iteration: (a) (a)

K=1 {27.6} = {22, 27, 35, 36} = 30

K=2 {9.5} = {9, 10, 11} = 10

K=3 {12.6} = {12, 15, 18} = 15

3rd Iteration:

K=1 {30} = {27, 35, 36} = 32.6

K=2 {10} = {9, 12, 11} = 10.5

K=3 {15} = {22, 15, 18} = 18.3

4th Iteration:

K=1 {32.6} = {27, 35, 36} = 32.6

K=2 {10.5} = {9, 12, 10, 11} = 10.5

K=3 {18.3} = {22, 15, 18} = 18.3

As the means in 3rd and 4th are same, which is stopping criteria for K-Means Algorithm, so clusters are -

$$C_1 = \{27, 35, 36\}$$

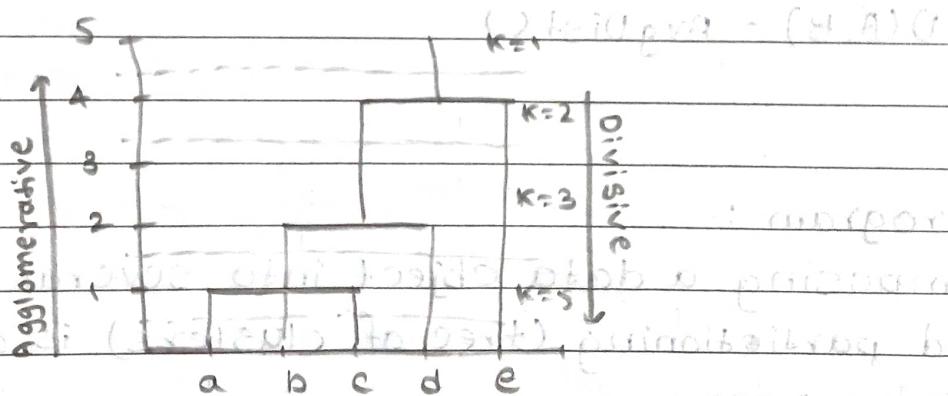
$$C_2 = \{9, 10, 11, 12\}$$

$$C_3 = \{15, 18, 22\}$$

Drawback of K-Means Algorithm :-

- ① The value of K should be given.
- ② It is sensitive to outliers.

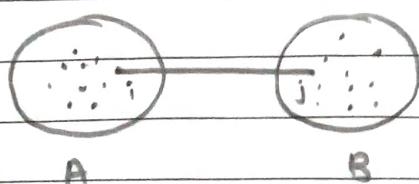
* Hierarchical Clustering :-



Dendrogram (to which the points are -)

* Agglomerative (class bottom up)

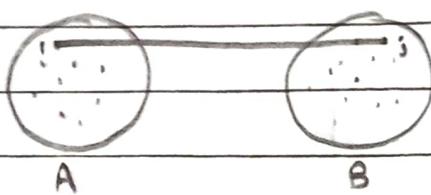
i) Single Linkage



$$D(A, B) = \min D(i, j)$$

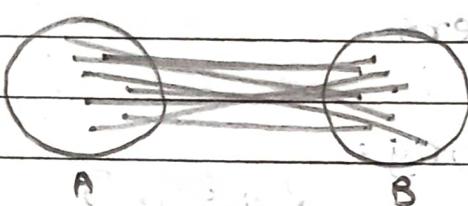
2) Complete Linkage

↳ Every object is connected with every other object.



$$D(A, B) = \text{MaxDist}(i, j)$$

3) Average Linkage



$$D(A, B) = \text{AvgDist}()$$

* Dendrogram :-

- Decomposing a data object into several levels of nested partitioning (tree of clusters) is called a Dendrogram.
- A clustering of data object is obtained by cutting the dendrogram, depending upon the value of k .
(at the derived level)
- Each connected component forms a cluster.

Q. Single Linkage

	1	2	3	4	5	Min Dist is 15.2
①	0					∴ Merge 1 & 2
2	2	0				
3	6	3	0			
4	10	9	7	0		
5	9	8	5	4	0	

S 15 P 01 A

	1,2	3	4	5	Min Dist is 15.3
②	0				∴ Merge 1,2 & 3
3	3	0			
4	9	7	0		
5	8	5	4	0	

S 15 P 01 A

	1,2,3	4	5	Min Dist is 4
③	0			∴ Merge 4 & 5
4	7	0		
5	5	4	0	

S 15 P 01 A

	1,2,3	4,5	
④	0		
4,5	5	0	

S 15 P 01 A

	1,2,3,4,5		
⑤	0		
1,2,3,4,5	0		

S 15 P 01 A

K=1

K=2

K=3

K=4

K=5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5

1 2 3 4 5</

SPECIALS STORE

Complete Linkage

Z A S S T

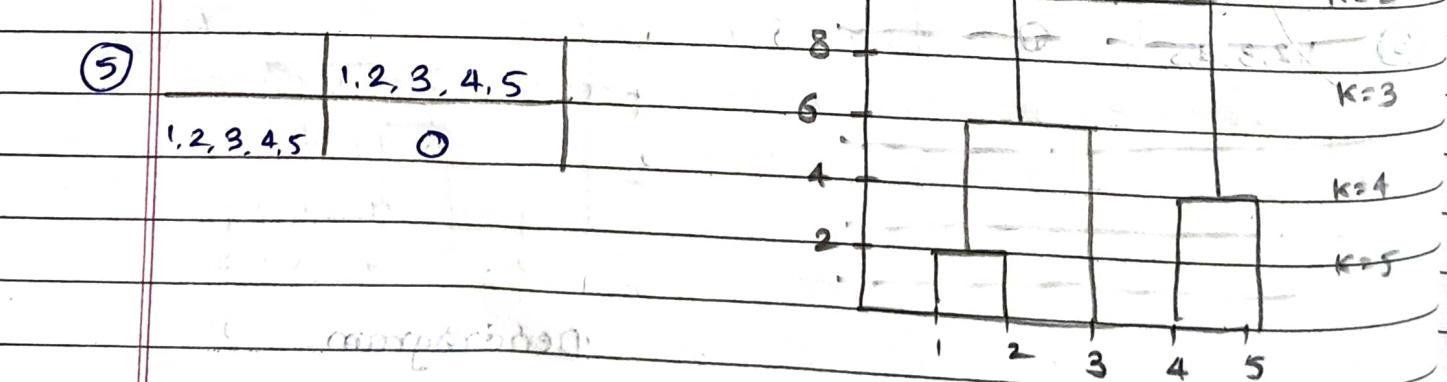
C D E

①	1	2	3	4	5	C	D	E
1	0					Min = 2	Z	A
2	2	0				Merge 1 & 2	A	
3	6	3	0			Z, A, S, H, C		
4	10	9	7	0				
5	98	8	5	4	0	Z, A, S, H, C		

②	1,2	3	4	5	C	D	E
1,2	0				Min = 4	Z	A
3	6	0			Merge 1 & 2	S	
4	10	7	0				
5	9	58	4	0	Z, A, S, H, C		

③	1,2	3	4,5	C	D	E
1,2	0			Min = 6	Z	A
3	6	0		Merge 1,2 & 3		
4,5	10	7	0	Z, A, S, H, C		

④	1,2,3	4,5	Merge 1,2,3,4,5	C	D	E
1,2,3	0			K=2		
4,5	10	0				



Average Linkage

①

	1	2	3	4	5	6	7	8	9	10
1	0									
2	2	0								
3	6	3	0					0	1	8
4	10	9	7	0	0		0	2	3	4
5	9	8	5	4	0		0	0	1	11

Min = 2

Merge 1,2

②

	1,2	3	4	5	6	7	8	9	10	
1,2	0									
3	4.5	0								
4	9.5	7	0							
5	8.5	5	4	0						

Min = 4

Merge 4,5

③

	1,2	3	4,5	6	7	8	9	10
1,2	0							
3	4.5	0						
4,5	9	6	0					

Min = 8.5

Merge 1,2 & 3

④

	1,2,3	4,5	6	7	8	9	10
1,2,3	0						
4,5	7.5	0					

Min = 7.5

Merge 1,2,3,4,5

⑤

	1,2,3,4,5	6	7	8	9	10
1,2,3,4,5	0					
6		0				

K=2

K=3

K=4

K=5



H.W.

standard approach

Q.

	1	2	3	4	5	6	7	8	9	10
1	0							0	1	
2	9	0					0	5	5	
3	3	7	0			0	8	3	2	
4	6	5	9	0	0	7	6	0	4	
5	11	10	2	8	0	2	8	0	2	

By using hierarchical clustering method, draw Dendrogram using single linkage and average linkage method.

*

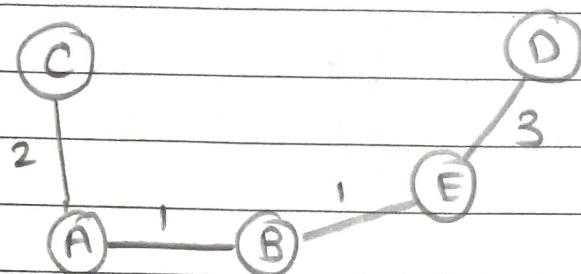
Divisive

	A	B	C	D	E	0	5	8	2
A	0					0	2	8	2
B	1	0				0	0	0	0
C	2	3	0						
D	5	4	6	0	2	4	8	8	2
E	4	1	2	3	0	0	0	8	2

- Hierarchical divisive algorithm uses minimum spanning Algo (Kruskal's principle) to create tree like structure which covers all nodes of connected undirected graph such that the sum of cost of edges is minimum.

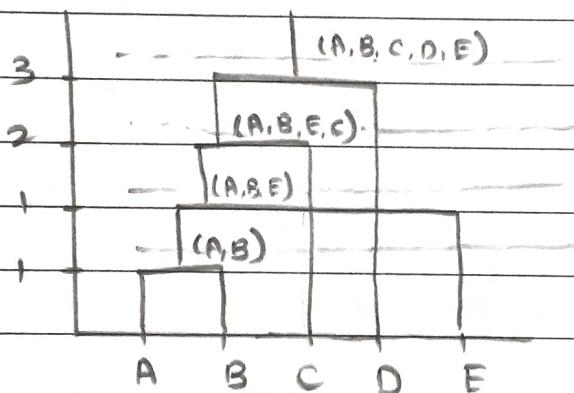
- MST never contains a cycle.
there can be many possible MSTs matrix, but they all have same minimum cost.

MST :-



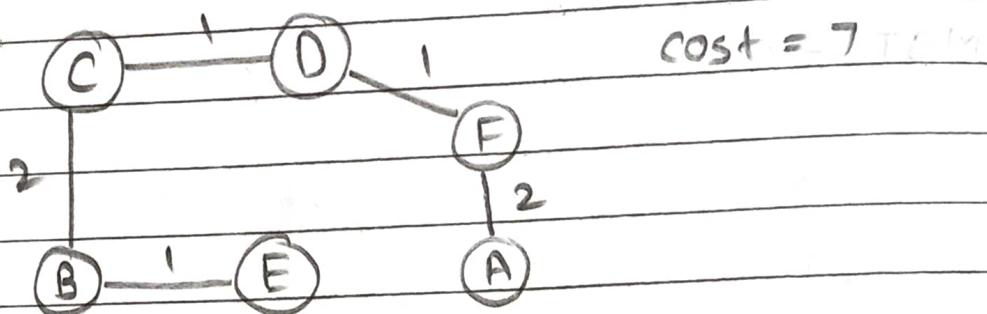
Chapman's Method

Dendrogram :-

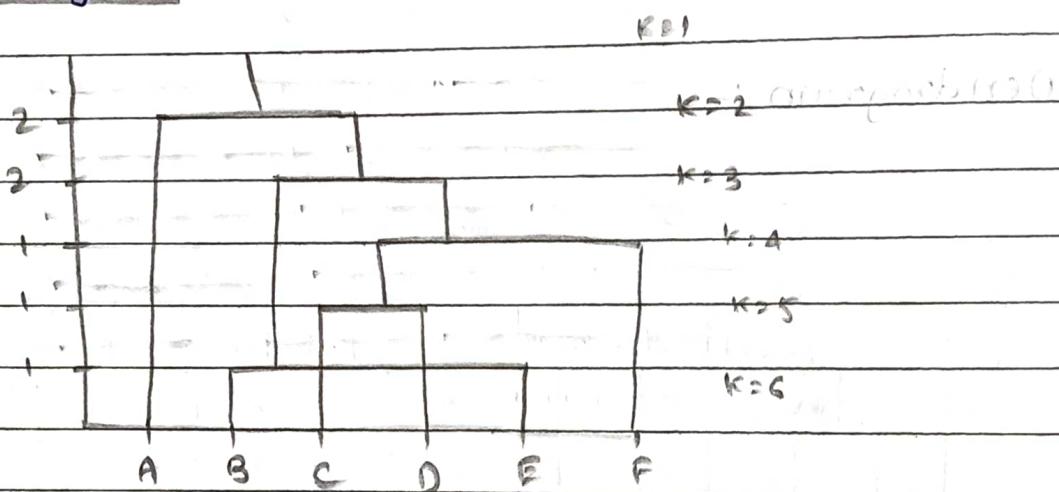


O.1	A	B	C	D	E	F
A	30					
B	3	0				
C	4	2	0			
D	7	3	1	0		
E	5	1	2	3	0	
F	2	6	7	1	4	0

MST :- Compare Ate Agglomerative & Divisive and create a Dendrogram using single & complete linkage method.



Dendrogram :-



* Association Mining :-

Frequent Patterns :-

- A pattern is a set of items that occur frequently in a dataset.
- A typical example of item set mining is market basket analysis (MBA).
- cross marketing, catalogue design, sale campaign analysis, web log analysis, DNA sequence analysis

* Apriori Algorithm

TWO Steps :

- 1) Join $L_1 \bowtie L_1$
- 2) Prune

TWO Lists :

- 1) Candidate List $\rightarrow C_1$
- 2) Item List $\rightarrow L_1$

Eg.-	Tid	Items		Itemset	Support
	10	I1 I3 I4	1st scan	I1 I3	2
	20	I2 I3 I5		I2 I3	3
	30	I1 I2 I3 I5	C1	I3 I5	3
	40	I2 I5		I4	1
				I5	3
		Sup Min = 2			

Item Set	Itemset		Itemset	support
I1	I1 I2	2nd scan	I1 I2	+
I2	I1 I3		I1 I3	2
I3	I1 I5	C2	I1 I5	1
I5	I2 I3		I2 I3	2
L1	I2 I5		I2 I5	3
	I3 I5		I3 I5	2

Item set	Itemset	3rd Scan	Itemset	Support
I1 I3	I1 I2 I5		I1 I2 I3	1
I2 I3	I1 I2 I3 I5	C3	I1 I2 I3 I5	1
I2 I5	I1 I3 I5		I1 I3 I5	1
I3 I5	I2 I3 I5		I2 I3 I5	2
L2	L2 M L2		C3	
			I2 I3 I5	
			Frequent Pattern	

Q.1	Tid	Items	SupMin = 2
1		I1 I2 I3	
2		I2 I4	
3		I2 I3	I3 I4 I5 I6
4		I1 I2 I4	
5		I1 I3	
6		I2 I3	
7		I1 I3	I2 I3 I4 I5 I6
8		I1 I2 I3 I5	
9		I1 I2 I3	

→ TCGA	Itemset	Support	Itemset	Support
	I1 S	6%	I1 I2	81 %
	I2 S	7%	I1 I3	81 %
	I3 I	6%	I2 I3	68 %
	I4 S	2%	I1 I4	04 %
	I5 S	2		2 = unique

TCGA	Itemset	Support	TCGA
1	S I1 I2	6%	1
2	S I1 I3	7%	1
3	S I1 I4	2%	1
4	S I2 I3	68 %	2
5	S I2 I4	2%	2
6	S I3 I4	04 %	3
7	I1 I2 I3	81 %	
8	I1 I2 I4	81 %	
9	I1 I3 I4	81 %	
10	I2 I3 I4	68 %	

TCGA	Itemset	Support	TCGA
1	S I1 I2 I3	81 %	1
2	S I1 I2 I4	81 %	2
3	S I1 I3 I4	81 %	3
4	S I2 I3 I4	68 %	

* Association Table Rules

Pattern : I₂ I₃ I₅

$$I_2 \Rightarrow I_3 I_5 = \frac{2}{3} = 66\%$$

$$I_3 \Rightarrow I_2 I_5 = \frac{2}{3} = 66\%$$

$$I_5 \Rightarrow I_3 I_2 = \frac{2}{3} = 66\%$$

$$I_2 \Rightarrow I_5 = \frac{3}{3} = 100\%$$

$$I_2 \Rightarrow I_3 = \frac{2}{3} = 66\%$$

$$I_3 \Rightarrow I_5 = \frac{2}{3} = 66\%$$

$$I_3 I_5 \Rightarrow I_2 = \frac{2}{2}$$

$$I_2 I_5 \Rightarrow I_3 = \frac{2}{3}$$

$$I_3 I_2 \Rightarrow I_5 = \frac{2}{3}$$

The rules whose confidence & support is greater than given threshold value, if it is called as strong association rule.

Drawback of Apriori :-

① TOO many scans.

Q. A database has 5 transactions. let min support be 60% & min confidence is 80%. Find all frequent item sets.

T100	'M', 'O', 'N', 'K', 'E', 'Y'
T200	'D', 'O', 'N', 'K', 'E', 'Y'
T300	'M', 'A', 'K', 'E'
T400	'M', 'U', 'C', 'K', 'Y'
T500	'C', 'O', 'O', 'K', 'I', 'E'

* FP growth

Minsupport = 2

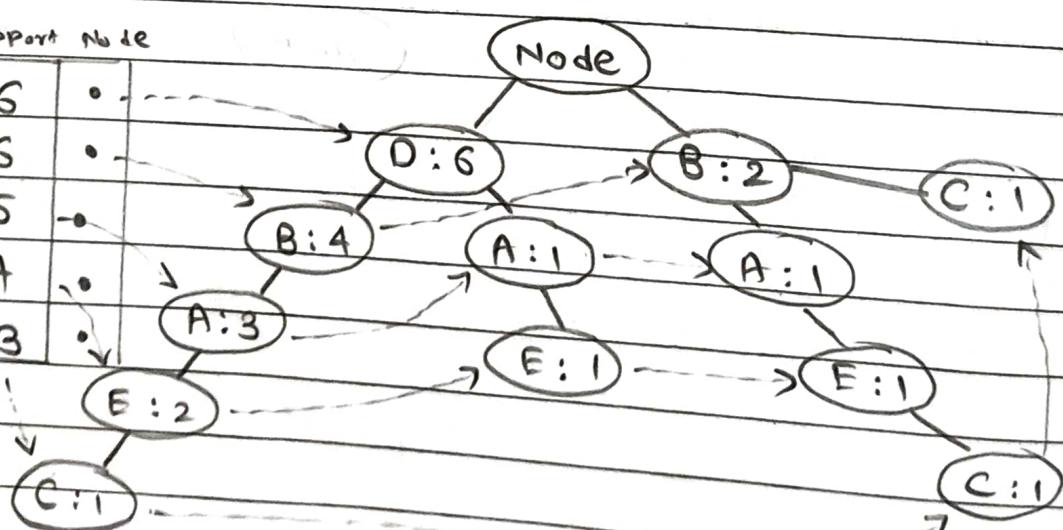
TID	Items	Itemset	SUPPORT
T1	E, A, D, B	A	5
T2	D, A, C, E, B	Scan → B	6
T3	C A B E	C	3
T4	B A D	D	6
T5	D	E	4
T6	D B		
T7	A D E		
T8	B C		

Decreasing
order

TID	Items	Itemset	SUPPORT
T1	D, B, A, E	D	6
T2	D, B, A, E, C	B	6
T3	B, A, E, C	Rearrange transactions ← A	5
T4	D, B, A	E	4
T5	D	C	3
T6	D, B		
T7	D, A, E		
T8	B, C		

Item Support Node

Item	Support	Node
D	6	•
B	6	•
A	5	•
E	4	•
C	3	•



FP growth Tree

Conversion of FP tree to Suffix form (LSD).

Condition	Frequent Base	FP
C	$\langle B:1 \rangle \langle EAB:1 \rangle \langle EABD:1 \rangle$	$\langle B:8 \rangle \langle EAB:2 \rangle \langle A:2 \rangle \langle E:2 \rangle \langle AB:2 \rangle \langle EA:2 \rangle$ $\langle E:2 \rangle \langle B:2 \rangle$
E	$\langle ABD:2 \rangle \langle AD:1 \rangle \langle AB:1 \rangle$	$\langle ABD:2 \rangle \langle A:4 \rangle \langle AB:3 \rangle \langle AD:3 \rangle \langle B:3 \rangle \langle D:3 \rangle$
A	$\langle BD:3 \rangle \langle D:1 \rangle \langle B:1 \rangle$	$\langle BD:3 \rangle \langle B:4 \rangle \langle D:4 \rangle$
B	$\langle D:4 \rangle \langle \rangle$	$\langle D:4 \rangle$
D		

Drawback :-

- Storage is required to store the tree.

Rules :-

$$C \Rightarrow B \quad E \Rightarrow ABD \quad A \Rightarrow BD \quad B \Rightarrow D$$

$$C \Rightarrow EAB \quad E \Rightarrow A \quad A \Rightarrow B$$

$$C \Rightarrow A \quad E \Rightarrow AB \quad A \Rightarrow D$$

$$C \Rightarrow AB \quad E \Rightarrow AD$$

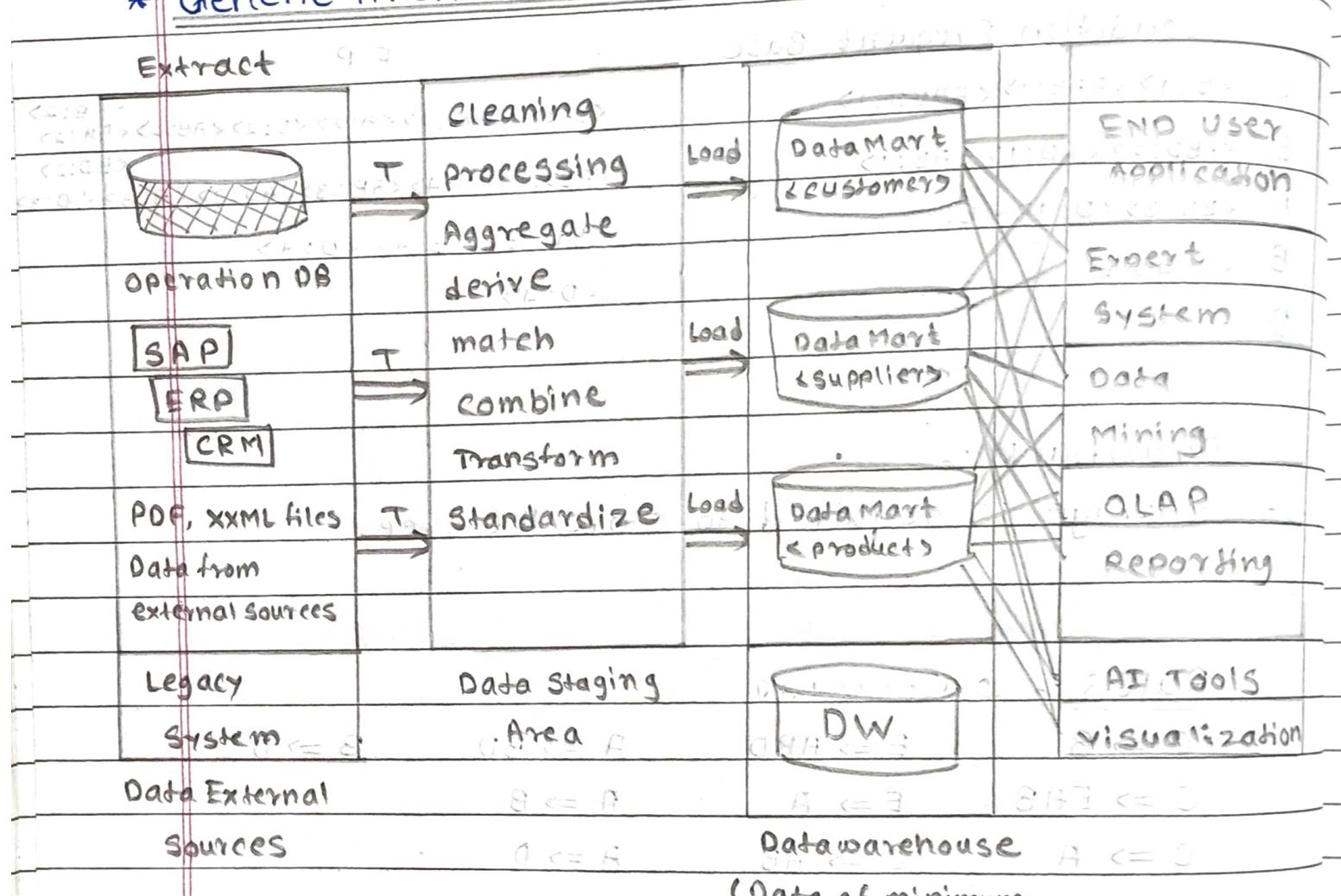
$$C \Rightarrow E \quad E \Rightarrow B$$

$$C \Rightarrow EB \quad E \Rightarrow BD$$

$$C \Rightarrow EA \quad E \Rightarrow D$$

Q.

* Generic Architecture of Data warehouse :-



Data warehouse :-

- ① Subject oriented
- ② Integrated
- ③ Time variant
- ④ Non Updatable (read only)

Data warehouse:-

A subj collection of data used in support of management decision making process.

Data warehousing:-

Process of constructing and using datawarehouse