

~~Ans~~

11/4/25

TID	Items	→ find all frequent itemsets using apriori algo. min-support = 2
1	{a, b, d, e}	
2	{b, c, d}	
3	{a, b, d, e}	
4	{a, c, d, e}	→ find all the association rules corresponding to frequent itemsets identified
5	{b, c, d, e}	in first part with given confidence % as <u>60%</u> .
6	{b, d, e}	
7	{c, d}	
8	{a, b, c}	
9	{a, d, e}	
10	{b, d}	→ Draw a lattice structure to represent all the association rule corresponding to frequent itemsets in ① and use confidence to identify all the association rules not satisfying condition confidence % > 60%.

1 itemset: itemset freq

{a} 5

{b} 6

{c} 5

{d} 4

{e} 6

2 itemset:	itemset freq	itemset freq	
{a, b}	3	{a, d}	4
{a, c}	2	{a, e}	4

$\{b, c\}$	3	$\{c, d\}$	4
$\{b, d\}$	5	$\{c, e\}$	2
$\{b, \cancel{e}\}$	3	$\{d, e\}$	6
$\{c, \cancel{e}\}$			

3 itemset:

$\{a, b, d\}$	2	$\{b, d, e\}$	3
$\{a, b, e\}$	2	$\{c, d, e\}$	2
$\{a, d, e\}$	4	$\{a, c, d\}$	1
$\{b, c, d\}$	2	$\{a, c, e\}$	1
$\{b, c, e\}$	1	$\{a, b, c\}$	1

4 itemset:

~~$\{a, b, c, d\}$ ,  $\{a, b, c, e\}$~~   $\rightarrow$  ~~abde~~  $\rightarrow$  ~~abde~~

~~$\{a, b, c, d\}$~~   $\rightarrow$  ~~abde~~

$\{a, b, c, e\} \rightarrow$  ~~abde~~

$\{a, b, d, e\} \rightarrow$  ~~abde~~  $\Rightarrow$  frequent set is

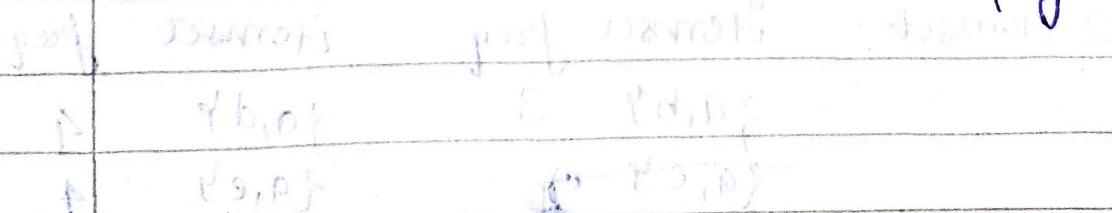
~~$\{a, d, c, e\}$~~   $\rightarrow$  ~~abde~~

~~$\{b, c, d, e\}$~~   $\rightarrow$  ~~abde~~

$\Rightarrow$  how to find association rules now?

We can partition abde in anyway.

one useful way to partition is using lattice structure. (as drawn on next page)

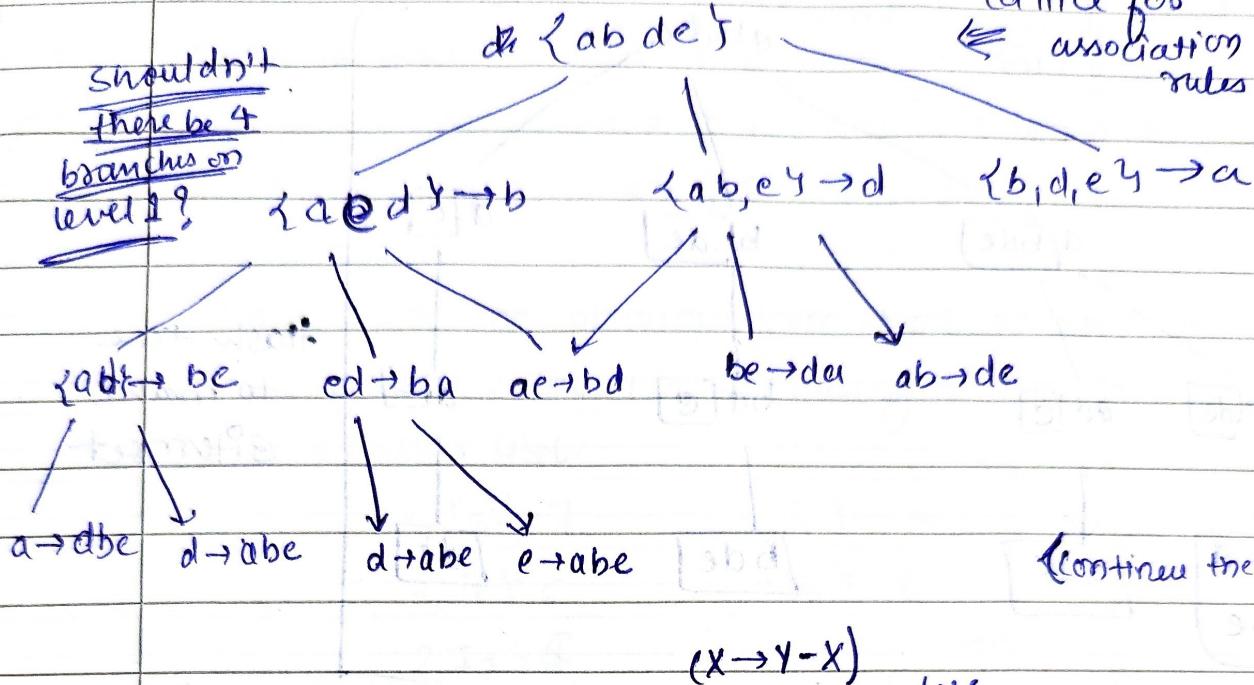


$$5C_3 = \frac{5!}{3!2!} = 10$$

$$5C_4 = \frac{5!}{4!1!} = 5$$

$$5C_3 = \frac{5!}{3!2!} = \frac{120}{2} = 60$$

shouldn't there be 4 branches on level 1?



$(X \rightarrow Y - X)$

- Now if confidence value of a rule is ~~less~~ than given threshold confidence value than any subset  $X'$  of  $X$  will also not satisfy the criteria

### Hash Base Support counting Method

- make hash tree for candidate itemset

ab

ac

ad

ae

bc

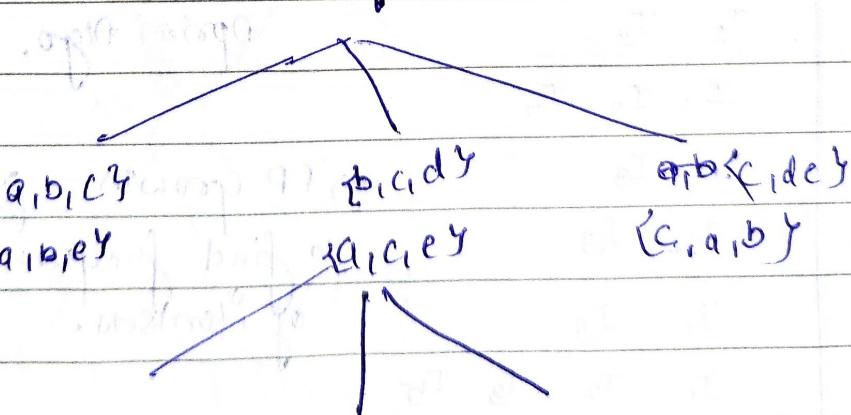
bd

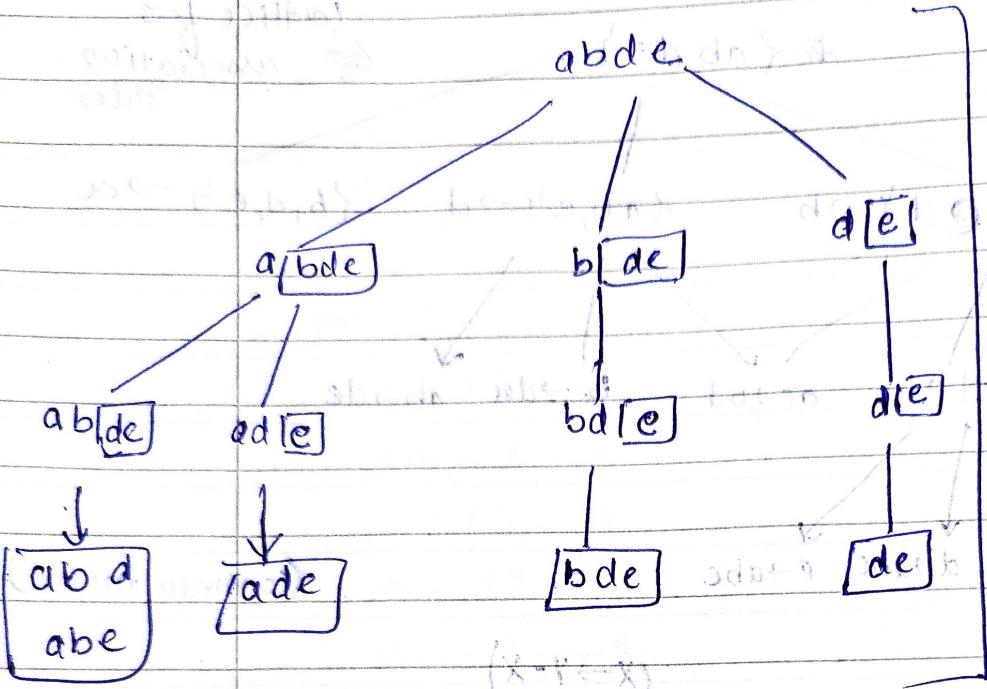
be

cd

ce

de





hash tree  
to make  
3-itemset

## FP Growth (frequent Pattern Growth)

18 | 4 | 25

with support = 2

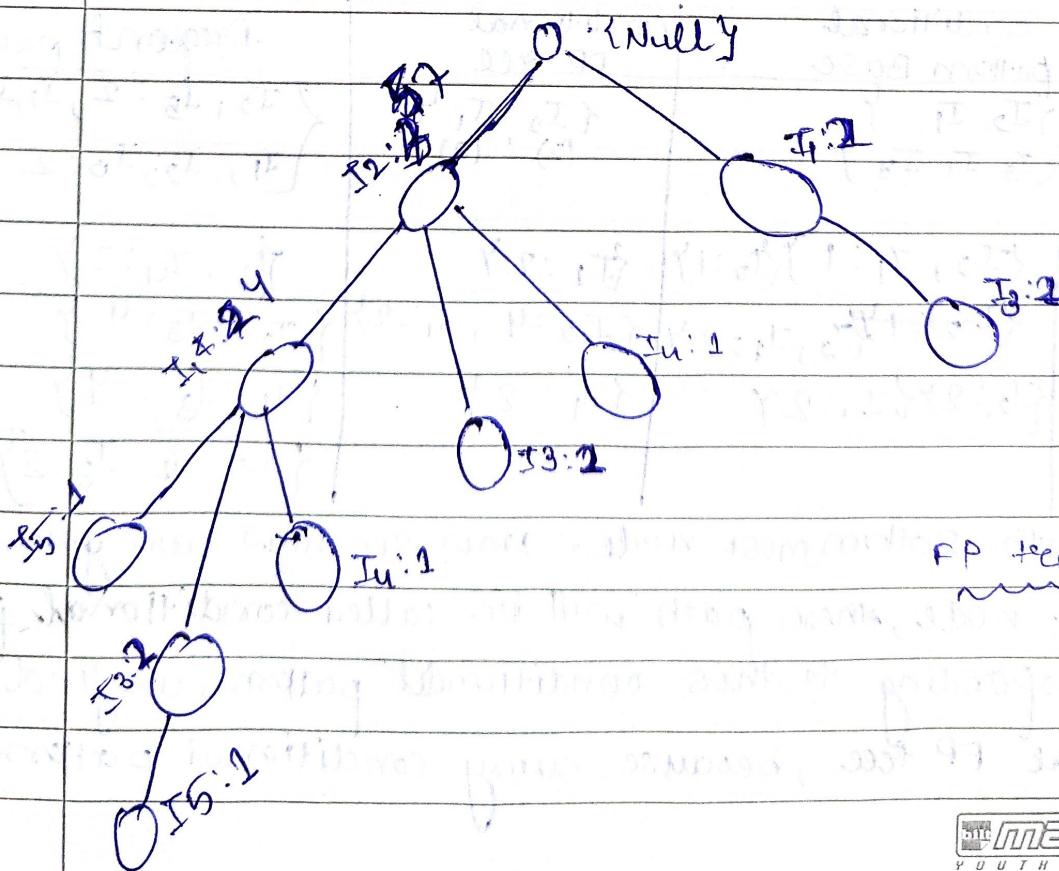
TID	items	→ need to scan table for counting support of items → drawback of Apriori algo.
T1	I <sub>1</sub> I <sub>2</sub> I <sub>5</sub>	
T2	I <sub>2</sub> I <sub>4</sub> I <sub>5</sub>	
T3	I <sub>2</sub> I <sub>3</sub>	
T4	I <sub>1</sub> I <sub>2</sub> I <sub>4</sub>	
T5	I <sub>1</sub> I <sub>3</sub>	→ FP growth → popular algo
T6	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>	to find frequent pattern of itemsets.
T7	I <sub>1</sub> I <sub>3</sub>	
T8	I <sub>1</sub> I <sub>2</sub> I <sub>3</sub> I <sub>5</sub>	
T9	I <sub>1</sub> I <sub>2</sub> I <sub>3</sub>	[Step 1] Scan the dataset &
T10	find support for all 1-itemset to find all frequent one itemset (Same as Apriori algo)	

	Support count	Support count
from step 1: $\{I_1\}$	6	9
$\{I_2\}$	7	2
$\{I_3\}$	6	2

**Step 2** Sort list obtained from step 1 as per their support values (in Decreasing Order). This is called Lorder list.

$\{I_2\} 7$        $\{I_1\} 6$        $\{I_3\} 2$   
 $\{I_2\} 7$        $\{I_1\} 6$        $\{I_3\} 2$   
 $\{I_3\} 6$

**Step 3** Prepare FP tree using Lorder list. consider null as root node & list all transactions as a branch in FP tree as per the order of the ~~Lorder~~ list.



Instead of seeing whole dataset as unit, divide it into subsets & find frequent itemset of each subset (main idea for FP tree)

⇒ To make the tree increment frequency of each item when a pattern can be traced on a branch else create a new branch.

⇒ following data structure is maintained:

items	support count	links
-------	---------------	-------

I1	7	
I2	6	
I3	6	
I4	2	
I5	2	

will include pointer  
to the respective item  
in tree (multiple pointers  
if item occurs in  
more than 1 node)

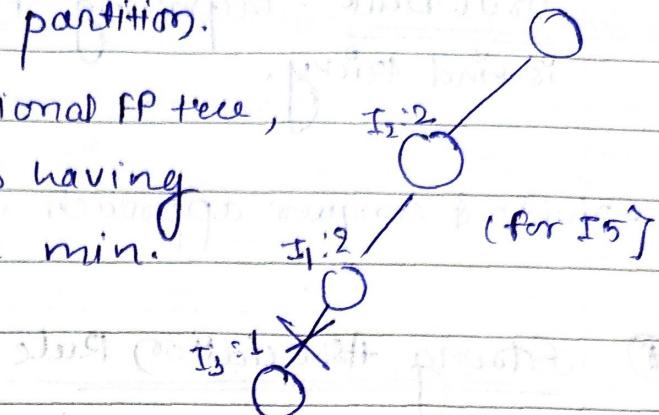
finding pattern using FP tree: start from bottom of the FP tree

items	conditional pattern base	conditional FP tree	frequent pattern
{I5}	{I2, I1}	{I2, I1, 4}	{I2, I5: 2, I1, I5: 2}
	{I2, I1, I3}	(2), (2)	{I1, I2, I5: 2}
{I4}	{I2, I1: 1} {I5: 1}	I2: 2	{I2, I4: 2}
{I3}	<del>{I2, I1: 1} {I5: 1}</del> {I2, I1: 2}	I2: 4, I1: 2	{I2, I3: 4}
	{I2: 2} {I1: 2}	I4: 2	{I1, I3: 4}
			{I2, I4, I3: 2}

Start with bottom most node. Now see how can you reach that node, these path will be called conditional pattern base. Corresponding to this conditional pattern, we find conditional FP tree, because every conditional pattern base

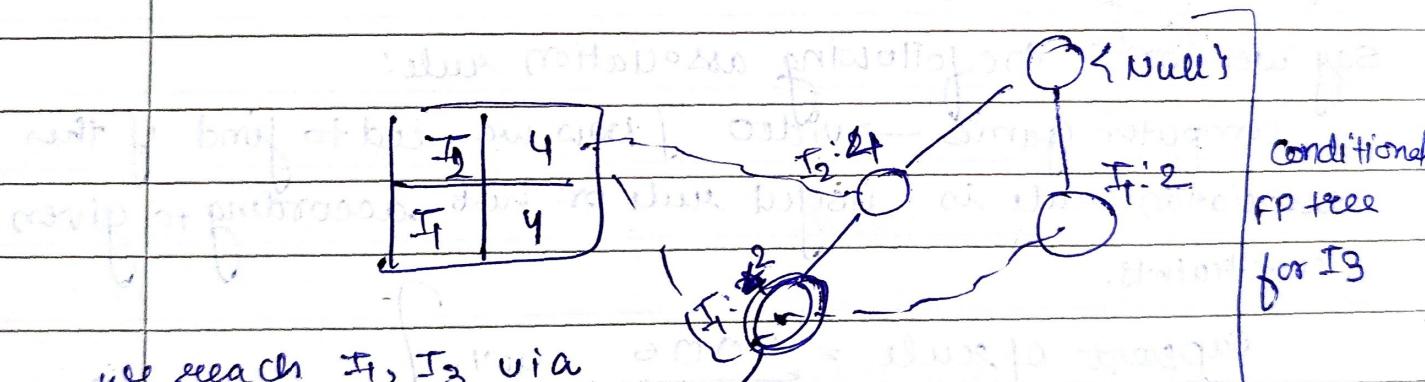
is considered as a partition.

Now form conditional FP tree, delete the nodes having support less than min. support count



concatenate the itemset with all combinations of conditional FP tree node.

Clearly not need to scan again for support values



this I1, so  
it will be 4 (adding 2 from this)

table entry for I1:

{I1}	{I2:2:4}	{I2:2:1} {I2:4}	{I2, I1:4}
------	----------	-----------------	------------

similarly find entries for all itemsets.

- The frequent itemsets obtained from the table are the same as ones found by Apriori algorithm.
- ⇒ This method lists all freq. itemsets in the dataset.

Drawback: Preparing FP tree with large dataset is kind tricky.

• Divide & conquer approach is used.

### (a) Strong Association Rule

$$N = 10000$$

$$\text{Computer Games} = 6000$$

$$\text{Video} = 7500 \quad \text{Computer Games \& Videos} = 4000$$

support=40% & confidence=66% ] Threshold values

Say we write the following association rule:

computer games  $\rightarrow$  video ] Now, we need to find if this association rule is useful rule or not according to given constraints.

$$\text{Support of rule} = \frac{4000}{10000} = 40\%$$

$$\text{confidence of Rule} = \frac{4000}{6000} = 66\%$$

} Satisfies the constraints, hence it is a useful rule

If videos are purchased alone  $\rightarrow$  their probability is 75% but if they are purchased with computer games the confidence is 66%  $\Rightarrow$  Negative impact of buying video with computer games. Hence, only support & confidence are not useful. We need association value to find items are positively associated or negatively

associated with each other. One measure to find this out is lift measure.

$$\text{lift}(A, B) = \frac{P(A \cup B)}{P(A) \cdot P(B)} \quad (A, B \rightarrow \text{two Itemset})$$

If  $\text{lift}(A, B) > 1 \Rightarrow$  Itemsets are Independent  
 $> 1 \Rightarrow$  Positive association  
 $< 1 \Rightarrow$  Negative association.

computing  $\text{lift}(\text{computer games}, \text{video}) = \frac{400}{1000}$

~~total number of document items~~ ~~lift measure~~  
 ~~$\frac{6000}{10000} \times \frac{7500}{10000}$~~   
~~= 40~~

~~$\frac{40}{60 \times 75} = 0.88$~~

~~$0.0000 = (a + b)$~~   
 ~~$0.0000 = (a + c)$~~   
 ~~$0.0000 = (b + d)$~~   
 ~~$0.0000 = (c + d)$~~   
This is less than 1, this means negative association b/w items.

Using  $\chi^2$  test we can find association here also.  
contingency table:

		Computer games		Total	
		Video	No Video	Computer games	No Computer games
Video	Video	4000 (4500)	3500 (3000)	3500 (3000)	2500
	No Video	2000 (1500)	500 (1000)	4000	10000
	Total	6000	8000	7500	12500

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

$$\begin{aligned}
 \chi^2 &= \frac{(4000 - 4500)^2}{4500} + \frac{(3500 - 3000)^2}{3000} + \frac{(500)^2}{1500} + \frac{(500)^2}{1000} \\
 &= \frac{(500)^2}{1500} \left( \frac{1}{3} + \frac{1}{2} + 1 \right) + \frac{500 \times 500}{1000} \\
 &= \frac{1}{3} \times 500 \left( \frac{2+3+6}{6} \right) + 500 = \cancel{\frac{650}{3}} = 216.667 \\
 &= \frac{500 \times 11}{6} + 50 = 555.55
 \end{aligned}$$

$\chi^2$  test considers hypothesis that itemsets are independent

Eg 2 N = 10000

$$n(CG) = 6000 \quad n(CG \& V) = 4000$$

$$n(V) = 7000 \quad n(V \& D) = 5500$$

$$n(D) = 8000 \quad n(CG \& D) = 5200$$

$$\text{min. support} = 40\%, \quad \text{min conf} = 66\%.$$

find DVD is likely to be bought with videos or computer games.

$$\text{lift}(D, V) = \frac{5500 / 10000}{\frac{1000}{10000} \times \frac{8000}{10000}} = \frac{5500 \times 10000}{1000 \times 8000} = \cancel{0.822}$$

$$\text{lift}(D, CG) = \frac{4000 \times 10000}{8000 \times 6000} = \cancel{0.833} < \frac{5200 \times 10000}{8000 \times 6000} = \cancel{0.833}$$

∴ DVD is more likely to be bought with videos.

# Cluster Analysis

15/4/25

clustering: Grouping of data objects  
→ unsupervised learning algorithm.

- ⇒ Classification methods (or supervised learning Algos) are called as learning by examples (because we have labeled data)
- ⇒ Clustering methods are called learning by observations.
- ⇒ Data objects in same cluster must be very similar to each other & data objects in different clusters must be very dissimilar.
- ⇒ Clustering is also used as data preprocessing method.  
(Summarization of data, Data compression,

vector quantization

→ Study

Good Quality Cluster → similarity of data objects within cluster should be very high.

cohesive → ~~Distilling~~ high intra-class similarity  
Distinctive → low inter-class similarity

What method to use to judge quality of a cluster?

## clustering

4 types of clustering algo :-

- (1) Partitioning Approach: Partitions are made into the dataset.
- (2) Hierarchical Approach (Need to have some hierarchical relationship b/w clusters).
- (3) Density Based Clustering: Based on density in particular area.
- (4)

### K-Mean clustering Algo (Partitioning Algo)

Partition dataset into  $K$  no. of clusters such that sum of squared distances is minimized.

$$E = \sum_{k=1}^K \sum_{p \in C_i} (p - c_i)^2 \quad \left\{ \begin{array}{l} c_i = \text{centroid of} \\ \text{cluster} \end{array} \right.$$

(22) 4 (25)

### Example for K Means Clustering Algo (K=2 given)

- 1) Select 2 random points data points.
- 2) Compute distance of rest of the points with previously chosen points. Cluster points acc. to distance from initial two points (keep it in cluster of closer point among the two chosen points). Following this step for all points generates two clusters.
- 3) Now compute mean value for each cluster.
- 4) These means will be new centroids & corresponding to these points you need to recluster the data points again (as told in step 2).  
Keep repeating steps 2 to 4 till
  - (i) No change in cluster points
  - OR (ii) No change in mean value.

- For measuring distance any measure can be considered like Euclidean Distance or Manhattan Distance etc.
- Data objects get clustered around the central data object. (This is the major intuition of K Means clustering algo).

Drawbacks: (i) Applicable only to objects in continuous n-dimensional space (what about nominal attributes, this algo won't work there).

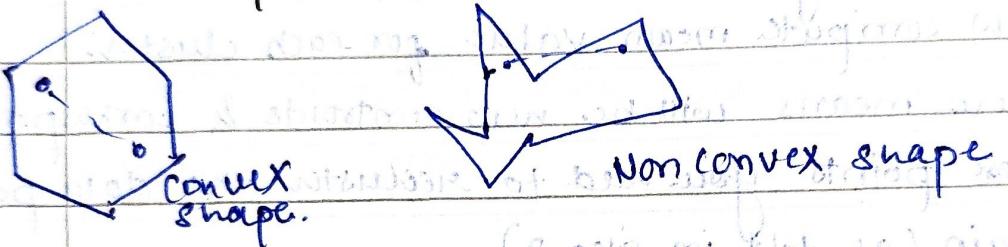
K-Modes: Works for nominal or categorical attributes. Instead of calculating mean, calculate mode value,

(rest of the algo is same as K Means Clustering)

2) Very sensitive towards outliers.

3) Not suitable to discover clusters with non convex shapes.

If we draw line to connect two points within a shape then that line must also lie within the shape  $\Rightarrow$  such shapes are called convex shapes. Otherwise called non convex shapes.



$$E = \sum_{i=1}^k \sum_{p \in C_i} (p - c_i)^2$$

Formula to compute quality of a cluster.  
(Within Cluster Variance)

If E is low  $\rightarrow$  Good quality cluster.

K Medoid: Instead of mean of cluster, we find the best representative element of the cluster and that element is considered medoid element.

Distance b/w central element & all other elements is min.

PAM method: uses medoid method

[Partition around Medoid]: starts from an initial set of medoids & iteratively replaces one of the medoids by one of the non medoids if it improves the total distance of resulting clustering.

## Example

$x_1$

$x_2$

→ Use Manhattan Distance as similarity measure

$O_1$

2

6

$O_2$

3

4

(A)

initial points

$O_3$

3

8

$O_4$

4

7

$O_5$

6

2

$O_6$

6

4

$O_7$

7

3

$O_8$

7

4

$O_9$

8

5

$O_{10}$

7

6

points

$O_1$

$O_3$

$O_4$

$O_5$

$O_6$

$O_7$

$O_9$

$O_{10}$

$d(A)$

3

4

4

5

3

3

1

1

2

2

$d(B)$

7

8

6

3

1

1

2

2

$$C_1 = \{O_2, O_1, O_3, O_4\} \quad C_2 = \{O_8, O_5, O_6, O_7, O_9, O_{10}\}$$

$$\text{Total Error} = E_1 + E_2$$

$$\text{Absolute Error : } E_1 = (O_1 - O_2) + (O_3 - O_2) + (O_4 - O_2)$$

$$= 3 + 4 + 4 = 11$$

$$E_2 = 3 + 1 + 1 + 2 + 2 = 9$$

$$E = E_1 + E_2 = 11 + 9 = \underline{\underline{20}}$$

Now acc. to PAM, let's replace  $O_8$  by (say)  $O_9$ .

New new cluster representative of  $C_2$  is  $O_9$ .

## ~~13/12/24~~ Distributed System Logs

→ Now recluster the data points with  $O_7$  &  $O_8$  as cluster points. We'll still get the same cluster structure.

points	$d(O_2)$	$d(O_7)$	$d(O_8)$	$d(O_9)$	$d(O_{10})$
$O_1$	3	8	8	8	8
$O_3$	4	9	9	9	9
$O_4$	4	7	7	7	7
$O_5$	5	2	2	2	2
$O_6$	3	2	2	2	2
$O_8$	4	1	1	1	1
$O_9$	6	3	3	3	3
$O_{10}$	6	3	3	3	3

Now recalculate error.

$$E_1 = 3 + 4 + 4 = 11$$

$$E = 22 \rightarrow \text{more}$$

$$E_2 = 2 + 2 + 1 + 3 + 3 = 11 \rightarrow \text{than old error}$$

If new error  $>$  old error  $\Rightarrow$  Don't accept  
else the cluster representative is selected.

Stopping condition:

→ PAM works well for small dataset but not for large one.

## Hierarchical Clustering

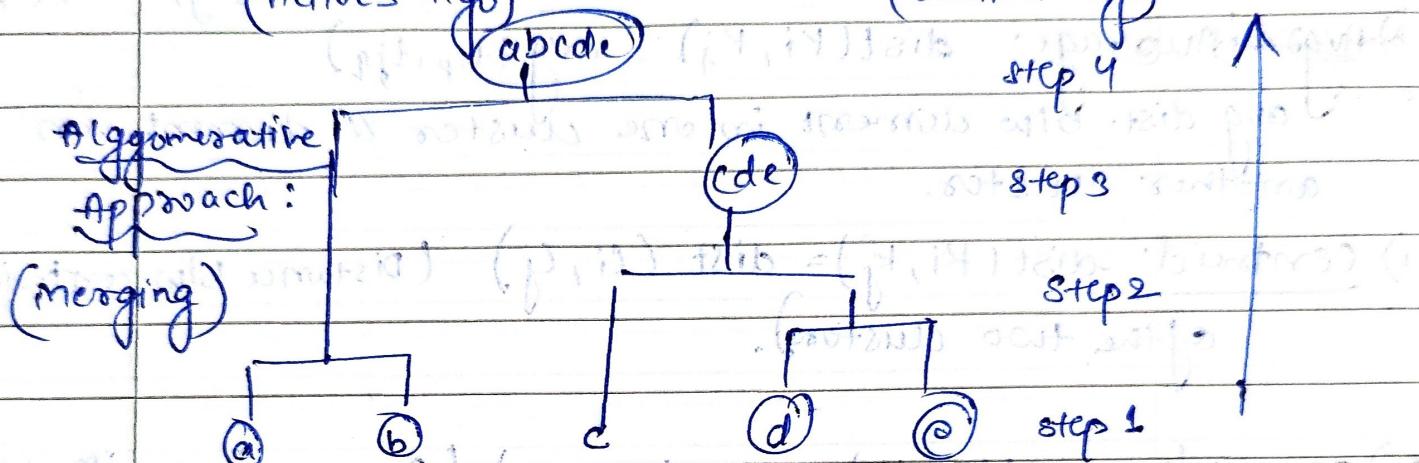
Two types

(1) Agglomerative (2) Divisive (splitting)

(Bottom to up approach)

(AGNES algo)

(Top to bottom approach)  
(DIANA algo)



Data points like a, b, c, d, e can either be single data points or cluster representative.

Note: question is should a be

→ finding best merge point is very important but very challenging also.

### AGNES (Agglomerative Nesting)

Merge one cluster with other if they have least dissimilarity (or max. similarity), (Various methods to measure similarity)

Dendrogram: Representation of how data objects in hierarchical clustering are associated.

Methods to find distance b/w clusters:

- 1) Single link: Min. dist. b/w an element in one cluster & an element in another cluster. i.e.  $\text{dist}(k_i, k_j) = \min(t_{ip}, t_{jq})$
- 2) Complete link: largest dist. b/w an element in one cluster & an element in another cluster i.e.  $\text{dist}(k_i, k_j) = \max(t_{ip}, t_{jq})$
- 3) Avg. Average:  $\text{dist}(k_i, k_j) = \text{avg.}(t_{ip}, t_{jq})$   
avg. dist. b/w element in one cluster & element in another cluster.
- 4) Centroid:  $\text{dist}(k_i, k_j) = \text{dist}(l_i, l_j)$  (Distance b/w centroids of the two clusters).
- 5) Medoid:  $\text{dist}(k_i, k_j) = \text{dist}(m_i, m_j)$  (Distance b/w medoids i.e. best representatives of the two clusters).