

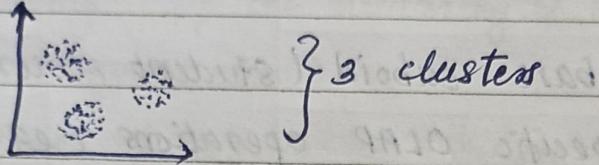
11.08.24

## MODULE 3 AGAIN

clustering

K-means clustering algorithm

divide data into groups - each group is a cluster.

spherical types of clusters - identified by K-means  
each cluster represented by a data object - clusterrepresentative / cluster ~~center~~ centre  
For K-means,

Representative - mean value of data objects in the cluster

K - total number of clusters

K = 2  $\rightarrow$  2 clusters; K = 3  $\rightarrow$  3 clusters

should be given as input to the algo.

K  $\rightarrow$  range : 1 to n-1

Here, 1 to 3.

let K be 2. Data objects: (1,1), (2,1), (4,3), (5,4)

Iteration 1:

	(1,1) Centroid 1 $\sqrt{(-1)^2 + (1-1)^2} = 0$	(2,1) Centroid 2 $\sqrt{(2-2)^2 + (1-1)^2} = 0$	cluster
(1,1)			1
(2,1)	$\sqrt{(2-1)^2 + (1-1)^2} = 1$	$\sqrt{(2-2)^2 + (1-1)^2} = 0$	2
(4,3)	$\sqrt{(4-1)^2 + (3-1)^2} = \sqrt{13} = 3.6$	$\sqrt{(4-2)^2 + (3-1)^2} = \sqrt{8} = 2.8$	2
(5,4)	$\sqrt{(5-1)^2 + (4-1)^2} = \sqrt{16+9} = 5$	$\sqrt{(5-2)^2 + (4-1)^2} = \sqrt{9+9} = 4.2$	2

Find proximity using Euclidean Distance Measure.

Inner one - clusters (in each row)

cluster 1 = { (1,1) }

cluster 2 = { (2,1), (4,3), (5,4) }

Iteration 2:		$(\frac{2+4+5}{3}, \frac{1+3+4}{3}) = (3.67, 2.67)$	
	$(1,1)$	Centroid 1	
$(1,1)$	0	$\sqrt{(1-3.6)^2 + (1-2.6)^2} = 3.05$	cluster 1
$(2,1)$	1	$\sqrt{(2-3.6)^2 + (1-2.6)^2} = 2.26$	1
$(4,3)$	3.6	$\sqrt{(4-3.6)^2 + (3-2.6)^2} = 0.57$	2
$(5,4)$	5	$\sqrt{(5-3.6)^2 + (4-2.6)^2} = 1.98$	2

since only one element in cluster 1 keep it as that. calculate mean of the other clusters and take that as centroid. If more than one element in one cluster, take the mean.

$$\text{cluster 1} = \{(1,1), (2,1)\}$$

$$\text{cluster 2} = \{(4,3), (5,4)\}$$

$$\text{Centroid 1} = (1,1) \quad \text{Centroid 2} (3.6, 2.6)$$

If in the next iteration, the <sup>cluster</sup> values change, we might have to continue further.

Iteration 3:

Centroid 1 (1.5, 1)		Centroid 2 (4.5, 3.5)	cluster
$(1,1)$	$\sqrt{(1-1.5)^2 + (1-1)^2} = 0.5$	$\sqrt{(4.5-1)^2 + (3.5-1)^2} = 4.3$	1
$(2,1)$	$\sqrt{(2-1.5)^2 + (1-1)^2} = 0.5$	$\sqrt{(4.5-2)^2 + (3.5-1)^2} = 3.5$	1
$(4,3)$	$\sqrt{(4-1.5)^2 + (3-1)^2} = 3.2$	$\sqrt{(4.5-4)^2 + (3.5-3)^2} = 0.7$	2
$(5,4)$	$\sqrt{(5-1.5)^2 + (4-1)^2} = 4.6$	$\sqrt{(4.5-5)^2 + (3.5-4)^2} = 0.7$	2

cluster 1 and cluster 2 remains the same.

$$\text{Centroid of cluster 1} = (1.5, 1)$$

$$\text{Centroid of cluster 2} = (4.5, 3.5)$$

Since Iteration 2 and 3 have the same output, we can stop.  
(otherwise, continue with next)

## Partitioning Clustering Algorithm

### Iterative Clustering Algorithm

In k-means, iterative relocation of data objects happens.

Disadvantage:

- the value of  $k$  is decided by user
- sensitive to outliers in dataset

Algorithm for k-means:

Repeat

(re)assign each data object to the cluster to which the object is most similar

Algorithm for k-means:

Input:  $K$  : no of clusters

$D$  : no of data objects

Output:  $K$  clusters

Arbitrarily choose  $K$  number of data objects from  $D$  as initial cluster centers.

Repeat

(re)assign each data object to the cluster to which the object is most similar

update the cluster means

until no change in the clusters

a. Data objects of dataset D:

$(1,2), (6,4), (4,9), (2,10), (2,5), (8,4), (5,8), (7,5)$

Find the clusters with  $k=3$ .

An. Iteration 1

	$(1,2)$ Centroid 1	$(5,8)$ Centroid 2	$(8,4)$ Centroid 3	Clusters
$(1,2)$	0	7.21	7.28	1
$(6,4)$	5.39	4.12	3	3
$(4,9)$	7.62	1.41	6.40	2
$(2,10)$	8.06	3.61	8.49	2
$(2,5)$	3.16	4.24	6.08	1
$(8,4)$	7.28	5	0	3
$(5,8)$	7.21	0	5	2
$(7,5)$	6.71	3.61	1.41	3

$$\text{cluster 1} = \{(1,2), (2,5)\}$$

$$\text{cluster 2} = \{(4,9), (2,10), (5,8)\}$$

$$\text{cluster 3} = \{(6,4), (8,4), (7,5)\}$$

Iteration 2

a. Dataset: 1, 2, 3, 20, 21, 22

Ans. Iteration 1

	Medoid 1 (1)	Medoid 2 (2)	Cost (min of other two rows)
1	0	1	0
2	1	0	0
3	2	1	1
20	19	18	18
21	20	19	19
22	21	20	20

Cost: 58

	Medoid [1, 3]	Medoid [1, 20]	Medoid [1, 21]	Medoid [1, 22]
1	2	19	20	21
2	1	0	19	20
3	0	0	18	19
20	17	17	1	2
21	18	18	0	1
22	19	19	2	0

new cost

55

5

4

5

new - old

-3

-53

-54

-53

replace 1 with 21

-54 is the suitable cost. Therefore replace 1 with 21.

	Medoid (2, 3)	Medoid (2, 20)	Medoid (2, 21)	Medoid (2, 22)
1	0	19	20	21
2	1	18	19	20
3	0	17	18	19
20	17	0	1	2
21	18	1	0	1
22	19	2	1	0

55  
-3 //

6  
-52 //

5  
-53 //

6  
-52 //

new medoid  
(2, 21)

disadvantage : inputs the value of  $k$ .

PAM is not suitable for high dimensional data.

### Hierarchical clustering

hierarchical decomposition of DB

2 types :

bottom up

top down

Iterative splitting of DB until some termination condition is satisfied.

### Density based clustering - Algorithm

[DBSCAN] → density Based Spatial clustering of Applications

advantage : can find arbitrarily shaped clusters

requires two parameters:

→ radius

1. eps - defines the neighbourhood around a data point.

2. Min Pts. - minimum number of neighbours.

minimum number of points in a circle with epsilon radius.

### Algorithm DBSCAN:

Input:

D : a data set containing

$\epsilon$  : the radius parameter, and

Min Pts : the neighbourhood density threshold

Output: A set of density based clusters

(1) mark all objects as unvisited.

(2) do

(3) randomly select an unvisited object  $p$ ;

(4) mark  $p$  as visited;

(5) if the  $\epsilon$ -neighborhood of  $p$  has at least  $M_{\text{pts}}$  objects

(6) create a new cluster  $C$ , and add  $p$  to  $C$ ;

(7) let  $N$  be the set of objects in the  $\epsilon$ -neighborhood of  $p$ ;

(8) for each point  $p'$  in  $N$

(9) if  $p'$  is unvisited

(10) mark  $p'$  as visited;

(11) if the  $\epsilon$ -neighborhood of  $p'$  has at least  $M_{\text{pts}}$  points, add those points to  $N$ ;

(12) if  $p'$  is not yet a member of any cluster, add  $p'$  to  $C$ ;

(13) end for

(14) output  $C$ ;

core object

iterative process

noise or outlier

$\epsilon$  must be chosen carefully.

In DBSCAN,

- core point

- datapoints

In DBSCAN Algorithm, three types of data objects.

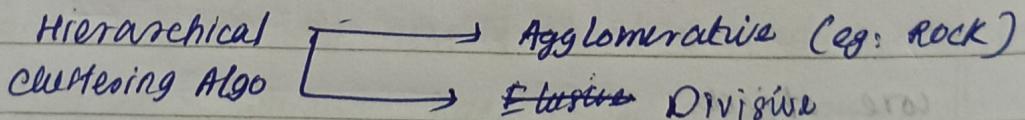
- Core Point
- Border Point
- noise / outliers

Direct Density Reachable (DDR) : A point  $p$  is DDR from point  $q$  if  $q$  is the core point and  $p$  is the neighbour of  $q$ .

Pensity Reachable (PR) : Two points are PR, if there is a chain of DDR points that link these two points.

### ROCK - Robust Clustering using links

hierarchical clustering algorithm that analyze the concept of links (the number of common neighbours among two objects) for data with categorical attributes.



Ways to find the proximity between clusters:

- single linkage (Euclidean distance between the closest pair of objects)
- complete linkage (distance between the two clusters = euclidean distance of the farthest pair of objects.)
- average linkage (distance is defined as the average of distances between all pairs of objects  $r$  and  $s$  belonging to different clusters.)

Find the proximity using single linkage.

(15, 11, 5, 10) Proximity = 0

(5, 1, 6, 2)

Proximity = 0

$$\sqrt{(5-5)^2}$$

## ROCK algorithm

1. obtaining a random sample of the data.
  2. Performing clustering on the sampled data using link agglomerative cluster approach where the goodness measure helps to identify the best pair of clusters to be merged.
  3. Assign the remaining data points of the disk to the generated clusters.

## Major definitions

- similarity function
  - Neighbors
  - Links
  - criterion function
  - goodness measure

26.08.24 MODULE 4

TID	List of item-IDs	Association Rule Mining
T100	I1, I2, I5	(I1, I2, I1, I1)
T200	I2, I4	(I2, I2, I2)
T300	I2, I3	(I2, I2, I2)
T400	I1, I2, I4	(I1, I1, I1)
T500	I1, I3	(I1, I1, I1)
T600	I2, I3	(I1, I1, I1)
T700	I1, I3	(I1, I1, I1)
T800	I1, I2, I3, I5	(I1, I1, I1)
T900	I1, I2, I3	(I1, I1, I1)

(Functional dependency vs.  
Association rules)  
↓  
changes  
not FDs

3 Types of Frequent Pattern → Frequent Itemset ~ set of items which are frequently purchased by the customer

- Market Basket Analysis
- Frequent Structured Pattern  
(subgraph that occurs frequently in a graph)
  - Frequent Sequential Pattern  
(list of items)

only interesting patterns give knowledge.

Measures for Rule Creation:

1. support ~> The support for  $X \rightarrow Y$  is the probability of both X and Y appearing together, that is  $P(X \cup Y)$
2. confidence

Absolute Support  $\frac{\text{norm}}{n}$  (total no of transactions in which item appears)

Relative support norm =  $\frac{n}{\text{total no of transactions}}$

$$\{ \text{Bread, chung} \} \cdot 2, 2/5 = 2, 0.4 \neq 40\%.$$

Confidence of  $X \rightarrow Y$  is the conditional probability of  $Y$  appearing given that  $X$  exists. It is written as  $P(Y|X)$  and read as P of  $Y$  given  $X$ .

$$\text{confidence}(A \Rightarrow B) = P(B|A)$$

$$= \frac{\text{support-count}(A \cup B)}{\text{support-count}(A)}$$

Two steps for mining association rules.

1) Find all frequent items.

2) Generate strong association rules from the frequent database.

1. Join

2. Pruning  $\rightsquigarrow$  based on a property called Apriori property.

Apriori property:

All subsets of a frequent itemset must be frequent.

1) Find frequent items  $L_1$ .

2)  $L_2 = L_1 \bowtie L_1$  join pruning

3)  $L_3 = L_2 \bowtie L_2$

4)  $L_4 = L_3 \bowtie L_3$

Find Association rules from the following data using Apriori algorithm with minimum support count is 2 and confidence. To %.

C1

itemset	sup. count	
I1	6	Compare support count with min. 8e.
I2	7	
I3	6	
I4	2	if support-count < min. se eliminate
I5	2	itemset sup. count

L1

itemset	sup. count	
I1	6	I1, I2 4
I2	7	$\Rightarrow$ I1, I3 4
I3	6	I1, I4 1
I4	2	I1, I5 2

L2

itemset	sup. count	
I1, I2	4	I2, I3 4
I1, I3	4	I2, I4 2
I1, I4	2	I2, I5 2
I1, I5	2	I3, I4 0
I2, I3	4	I3, I5 1
I2, I4	2	I4, I5 0

01-04-'24

C3 itemset	sup. count
I1, I2, I3	2
I1, I3, I5	1
I1, I2, I5	2
I2, I4, I5	0
I2, I3, I4	0
I2, I3, I5	1

Transactions = L3

I1, I4, I5	I1, I2, I3	2
I1, I2, I5		2

i.e.,	C <sub>3</sub>	itemset	sup count	L <sub>3</sub>
	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>	2		itemset 8-C.
	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>	2		I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> 2
	I <sub>1</sub> , I <sub>3</sub> , I <sub>5</sub>	1 + prune		I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub> 2
	I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub>	0 + prune		
	I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>	1 + prune		
	I <sub>2</sub> , I <sub>4</sub> , I <sub>5</sub>	0 + prune		subset not in L <sub>2</sub> .

Form the association rules

subsets of the frequent set  $\Rightarrow$  l-s (remaining set)

$\rightarrow$  > 70%  $\approx$  strong association rule

$$I_1 \Rightarrow \{I_2, I_3\} \quad \text{confidence} = \frac{\text{sup}(I_1, I_2, I_3)}{\text{sup}(I_1)} = \frac{2}{6} = 33.3\%$$

$$I_2 \Rightarrow \{I_1, I_3\} \quad \text{confidence} = \frac{2}{7} = 28.6\%$$

$$I_3 \Rightarrow \{I_1, I_2\} \quad \text{confidence} = \frac{2}{6} = 33.3\%$$

$$\{I_1, I_2\} \Rightarrow I_3 \quad \text{confidence} = \frac{2}{4} = 50\%$$

$$\{I_1, I_3\} \Rightarrow I_2 \quad \text{confidence} = \frac{2}{4} = 50\%$$

$$\{I_2, I_3\} \Rightarrow I_1 \quad \text{confidence} = \frac{2}{4} = 50\%$$

$$I_1 \Rightarrow \{I_2, I_5\} \quad \text{confidence} = \frac{2}{6} = 33.3\%$$

$$I_2 \Rightarrow \{I_1, I_5\} \quad \text{confidence} = \frac{2}{7} = 28.6\%$$

$$I_5 \Rightarrow \{I_1, I_2\} \quad \text{confidence} = \frac{2}{2} = 100\% \rightarrow \text{strong}$$

$$\{I_1, I_2\} \Rightarrow I_5 \quad \text{confidence} = \frac{2}{4} = 50\%$$

$$\{I_1, I_5\} \Rightarrow I_2 \quad \text{confidence} = \frac{2}{2} = 100\% \rightarrow$$

$$\{I_2, I_5\} \Rightarrow I_1 \quad \text{confidence} = \frac{2}{2} = 100\% \rightarrow$$

## Apriori Example

Transaction ID

	Items
T1	(H) Hot Dogs, (B) Buns, (K) Ketchup
T2	(H) Hot Dogs, (B) Buns
T3	(H) Hot Dogs, (C) Coke, (Ch) Chips
T4	(Ch) Chips, (Co) Coke
T5	(Ch) Chips, (K) Ketchup
T6	(H) Hot Dogs, (C) Coke, (Ch) Chips

$$\frac{33.33 \times 1}{100} = 2$$

minimum support threshold : 33.33 Minimum confidence threshold : 60

AN:

Transaction ID

Hot Dogs T1 B, H, K

Buns T2 B, H

Ketchup T3 Ch, Co, H

Coke T4 Ch, Co

Chips T5 Ch, K

T6 Ch, Co, H

C1 : L1 :  $\{B\} \subseteq \{B, H, K\}$

B 2 B 2

Ch 4 Ch 4  $\{B, C\} \subseteq \{B, H, K\}$

Co 3 Co 3  $\{B, C\} \subseteq \{B, H, K\}$

H 4 H 4  $\{B, H\} \subseteq \{B, H, K\}$

K 2 K 2  $\{B, K\} \subseteq \{B, H, K\}$

C2  $\{B, C\} 0 \quad \{B, H\} 0 \quad \{B, K\} 1 \quad \{C, H\} 2 \quad \{C, K\} 1$

$\{B, C\} 0 \quad \{B, H\} 2 \quad \{B, K\} 1 \quad \{C, H\} 2 \quad \{C, K\} 1$

$\{B, H\} 2 \quad \{B, K\} 1 \quad \{C, H\} 3 \quad \{C, K\} 0$

$\{B, K\} 1 \quad \{C, H\} 2 \quad \{C, K\} 0$

$$\begin{matrix} & \\ & \downarrow \\ & 1 \end{matrix}$$

L2 :

B, H 2

Ch, Co 3

Ch, H 2

Co, H 2

C3 :

Ch, Co, H 2

L3 :

Ch, Co, H 2

$$Ch \Rightarrow \{Co, H\}$$

$$Co \Rightarrow \{Ch, H\}$$

$$H \Rightarrow \{Ch, Co\}$$

$$\{Ch, Co\} \Rightarrow H$$

$$\{Co, H\} \Rightarrow Ch$$

$$\{Ch, H\} \Rightarrow Co$$

$$\text{confidence} = \frac{2}{4} = 50\%$$

$$\text{confidence} = \frac{2}{3} = 66.66\% \text{ and } \left. \begin{array}{l} \\ \end{array} \right\} \text{strong}$$

$$\text{confidence} = \frac{2}{4} = 50\%$$

$$\text{confidence} = \frac{2}{3} = 66.66\% \rightsquigarrow$$

$$\text{confidence} = \frac{2}{2} = 100\% \rightsquigarrow$$

$$\text{confidence} = \frac{2}{2} = 100\% \rightsquigarrow$$

3.4.24

Pincer-Search Algorithm

Improvement over Apriori algorithm.

Algorithm:
 $L_0 = \emptyset; k=1; C_1 = \{\{i\} | i \in I\}; S_0 = \emptyset;$ 
 $MFCS := \{\{1, 2, \dots, n\}\}, MFS = \emptyset;$ 
do until  $C_k = \emptyset$  and  $S_{k-1} = \emptyset$ read database and count supports for  $C_k$  and MFCS;
 $MFS := MFS \cup \{\text{frequent itemsets in MFCS}\};$ 
 $S_k := \{\text{infrequent itemsets in } C_k\};$ 
call MFCS-gen algorithm if  $S_k \neq \emptyset$ ;

call MFS-pruning procedure;

generate candidates  $C_{k+1}$  from  $C_k$ ; (similar to a priori's generate & prune)If any frequent itemset in  $C_k$  is removed in MFS-pruning procedurecall the recovery procedure to recover candidate  $\overset{to}{C}_{k+1}$ call MFCS prune procedure to prune candidates in  $C_{k+1}$  $k = k + 1;$ 

return MFS;

Q abc def

abcg

minimum support = 2

abdh

bcde

abc

$S$  represents the infrequent itemsets  
 $L$ -frequent itemsets from  
C.  
candidate itemsets

MFCs  $\rightarrow$  Maximal Frequent candidate <sup>Item</sup> Set

$$\text{MFCs} = \{a, b, c, d, e, f, g, h\}$$

MFS  $\rightsquigarrow$  Maximal Frequent Set

contains the frequent itemsets from MFCs

1) Form C1: Itemset sup-count

a	4	MFCs sup-count = 0
b	5	not frequent
c	4	
d	3	MFS = $\emptyset$
e	2	$S_k = \{a, b, d, g\} - \{h\}\}$
f	1	
g	1	
h	1	

MFCs-gen

for all itemsets  $s \in S_k$

for all itemsets  $m \in \text{MFCs}$

- if  $s$  is a subset of  $m$

$$\text{MFCs} = \text{MFCs} \setminus \{m\};$$

for all items  $e \in \text{domain } s$

if  $m \setminus \{e\}$  is not a subset of any itemset in  $S_k$

$$\text{MFCs} = \text{MFCs} \cup \{m \setminus \{e\}\};$$

return MFCs

L1	Itemset	Sup-Count
a		4
b		5
c		4
d		3
e		2

considering {f,g}

$$MFCS = \{a, b, c, d, e, g, h\}$$

considering {g,h}

$$MFCS = \{a, b, c, d, e, h\}$$

considering {b,f}

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

### MFS - Prune

for all itemsets c in  $C_k$

If c is a subset of any itemset in the current NES  
delete c from  $C_k$ ;

### MFCS - Prune

for all itemsets c in  $C_{k+1}$

If c is not a subset of any itemset in the current MFCS  
delete c from  $C_{k+1}$ ;

C <sub>2</sub>	itemset	Sup-count
----------------	---------	-----------

$\{a, b\}$	$4 / 20 = 0.20$
$\{a, c\}$	$3 / 20 = 0.15$
$\{a, d\}$	$2 / 20 = 0.10$
$\{a, e\}$	$1 / 20 = 0.05$

b, c 4

b, d 3

b, e 2

c, d 2

c, e 2

d, e 2

Recovery

for all itemsets  $I \in C_k$

for all itemsets  $m \in MFS$

if the first  $k-1$  items in  $I$  are also in  $m$

if suppose  $m.item = I.item_{k-1}$

for  $i$  from  $j+1$  to  $|m|$

$C_{k+1} = C_{k+1} \cup \{I.item_1, I.item_2, \dots, I.item_{k-1}, m.item_i\}$

$m.item_i\}$

8.04.24  $k=2$

$MFC8 = \{a, b, c, d, e\}$

$MFS = \phi \cup \phi = \phi$

$S_2 = \{\{a, e\}\}$

$L_2 = \{\{a, b\}, \{a, c\}, \{a, d\}, \{b, c\}, \{b, d\}, \{b, e\}, \{c, d\}, \{c, e\}, \{d, e\}\}$

call  $MFC8$ -gen algo:

$S \in \{\{a, e\}\}$

$m \in \{\{a, b, c, d, e\}\}$

$S \subseteq m$

$MFC_S = \emptyset$

$e \in \{\{a, e\}\}$

$m \setminus \{e\} \Rightarrow m = \{\{a, b, c, d\}, \{a, b, c, d\}\}$

$MFC_8 = \emptyset \cup \{\{b, c, d, e\}, \{a, b, c, d\}\}$

$MFC_S - \{a, b, c, d, e\}$

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

$MFS = \emptyset \cup \emptyset = \emptyset$

$MFC_8 = \{\{b, c, d, e\}, \{a, b, c, d\}\}$

$MFS - \text{prune} \Rightarrow \text{no action}$

$C_3 = \{\{a, b, c\}, \{a, b, d\}, \{a, c, d\}, \{b, c, d\}, \{b, c, e\}, \{b, d, e\}, \{c, d, e\}\}$

call recovery  $\Rightarrow$  no action

call  $MFC_8 - \text{prune} \Rightarrow \text{no action}$

$k=8$

$C_3:$  idemset sup-count

a, b, c 3

a, b, d 2

a, c, d 1

b, c, d 2

b, c, e 2

b, d, e 2

c, d, e 2

$$MFC8 = \{\{b, c, d, e\} = 2, \{a, b, c, d\} = 1\}$$

$$MFS = \phi \cup \{b, c, d, e\} = \{\{b, c, d, e\}\}$$

$$S_3 = \{\{a, c, d\}\}$$

$$L_3 = \{\{a, b, c\}, \{a, b, d\}, \{b, c, d\}, \{b, c, e\}, \{b, d, e\}, \{c, d, e\}\}$$

$$MFC8\text{-gen } S \in \{\{a, c, d\}\}$$

$$m \in \{\{b, c, d, e\}, \{a, b, \overline{c, d}\}\}$$

$$S \subseteq m$$

$$MFC8 = \{b, c, d, e\}$$

$$c \in S_3$$

$$m \setminus \{e\} = \{b, c, d\}, \{a, b, d\}, \{a, b, c\} \notin MFC8$$

$$MFC8 = \{\{b, c, d, e\}, \{b, c, d\} \setminus \{a, b, d\}, \{a, b, c\}\}$$

$$MFS = \text{MFS}_1 - MFS = \phi$$

MFS prune

$$L_3 = \phi$$

all elements are deleted from  $L_3$

generate  $C_4 \Rightarrow$  no action as  $L_3 = \phi$

call recovery  $\Rightarrow$  no action

call MFC8-prune  $\Rightarrow$  no action

stop.

Find the frequent itemsets using Pinser Search

$$\text{Sup-count} = 2$$

TID	list
T100	11, 12, 15
T200	12, 14
T300	12, 13
T400	11, 12, 14
T500	11, 13
T600	12, 13
T700	11, 13
T800	11, 12, 13, 15
T900	11, 12, 13

Compare Aprori and FP Growth (FBH vs A)

- 1) Pattern generation  
constructs  
FP tree, pairing items into singletons, pairs, triplets
- 2) Candidate generation  
no, yes
- 3) Memory usage  
compact version of database is saved, candidate combinations are saved in memory.
- 4) Process

## Partition Algorithm

modification

2 Phases :

Phase 1 - entire database / dataset is partitioned into disjoint sets - frequent itemsets (<sup>local large itemsets</sup>)  
 Phase 2 - merged to candidate set

Algorithm:

$C_k^P$  set of local candidate k-itemsets in partition p

$L_k^P$  set of local large k-itemsets in partition p

$L^P$  set of all local large itemsets in partition p

$C_k^G$  set of global candidate k-itemsets

$C^G$  set of all global candidate itemsets

$L_k^G$  set of global large k-itemsets.

- 1)  $P = \text{partition-database}(D)$
- 2)  $n = \text{number of partitions}$
- 3) for  $i = 1$  to  $n$  begin
- 4)     read-in-partition ( $p_i \in P$ )
- 5)      $L' = \text{gen-large-itemsets}(p_i)$
- 6) end
- 7) for ( $i=2$ ;  $L_i^3 \neq \emptyset$ ,  $j = 1, 2, \dots, n$ ;  $i++$ ) do
- 8)      $C_i^G = \bigcup_{j=1,2,\dots,n} L_j^3$  // Merge Phase
- 9) for ( $i=1$  to  $n$ ) begin // Phase II
- 10)     read-in-partition ( $p_i \in P$ )
- 11)     for all candidates  $c \in C^G$   $\text{gen-count}(c, p_i)$
- 12) end

} Partition algorithm

- Initially the database D is logically partitioned into n partitions.
- Phase I of the algorithm takes n iterations.
- During iteration i only partition  $p_i$  is considered.
- The function gen-large-itemsets takes a partition  $p_i$  as input and generates local large itemsets of all lengths,  $L_1^i, L_2^i, \dots, L_k^i$  as the output.
- In the merge phase, the local large itemsets of same lengths from all n partitions are combined to generate the global candidate itemsets.
- In phase II, the algorithm sets up counters for each global candidate itemsets, and count their support for the entire database and generates the global large itemsets.

### Problem

<u>Transaction</u>	<u>Itemset</u>
T1	12, 15
T2	12, 14
T3	14, 15
T4	12, 13
T5	15
T6	12, 13, 14

Partitions:  $\{T1, T2\}, \{T3, T4\}, \{T5, T6\}$

First Scan min-sup-count = 1

$\{T_1, T_2\}$   $I_1 - 1, I_2 - 1, I_4 - 1, I_5 - 1$

$\{I_1, I_5\} = 1, \{I_2, I_4\} = 1$

$\{T_3, T_4\}$   $I_4 - 1, I_5 - 1, I_2 - 1, I_3 - 1$

$\{I_4, I_5\} = 1, \{I_2, I_3\} = 1$

$\{T_5, T_6\}$   $I_2 - 1, I_3 - 1, I_4 - 1, I_5 - 1, \{I_2, I_3, I_4\} = 1$

$\{I_2, I_3\} = 1, \{I_3, I_4\} = 1, \{I_2, I_4\} = 1$

Second Scan min-sup-count = 2

$I_1 - 1^*$   $I_3 - 2$   $I_5 - 2$

$I_2 - 3$   $I_4 - 3$

$\{I_1, I_5\} = 1^*$   $\{I_2, I_4\} = 2$   $\{I_4, I_5\} = 1^*$   
 $\{I_2, I_3\} = 2$   $\{I_3 - I_4\} = 1^*$   $\{I_2, I_3, I_4\} = 1^*$

Shortlisted :  $I_2, I_3, I_4, I_5, \{I_2, I_3\}, \{I_2, I_4\}$

### Dynamic Itemset Counting Algorithm

Alternative to Apriori Itemset Generation (incremental approach)

Itemsets are dynamically added and deleted as transactions are read.

Solid box  $\square$   $\begin{cases} \text{confirmed frequent itemset (finished counting, } > \text{minsup)} \\ \text{infrequent itemset (fc., } < \text{minsup)} \end{cases}$

Solid circle  $\circ$   $\begin{cases} \text{confirmed infrequent itemset (fc., } < \text{minsup)} \\ \text{suspected frequent itemset (still counting, } > \text{minsup)} \end{cases}$

Dashed box  $\boxed{\quad}$   $\begin{cases} \text{suspected frequent itemset (still counting, } > \text{minsup)} \\ \text{infrequent itemset (s.c., } < \text{minsup)} \end{cases}$

## 2. DIC Algorithm

1. Mark the empty itemset with a solid square. Mark all the 1-itemsets with dashed circles. Leave all other itemsets unmarked.
2. While any dashed itemsets remain:
  1. Read M transactions (If we reach the end of the transaction file, continue from the beginning). For each transaction, increment the respective counters for the itemsets that appear in the transaction and are marked with dashes.
  2. If a dashed circle's count exceeds minsupp, turn it into a dashed square. If any immediate superset of it has all of its subsets as solid or dashed squares, add new counters for it and make it a dashed circle.

### Example :

$$\text{minsupp} = 25\% \quad (25\% \times 4 = 1 \rightarrow \text{minsupp})$$

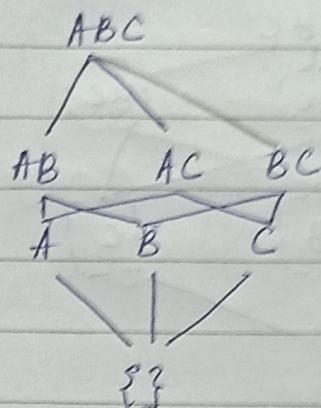
$$M = 2$$

TD	A	B	C	
T1	1	1	0	AB
T2	1	0	0	A
T3	0	1	1	BC
T4	0	0	0	

Transaction Database

Itemset lattices: an itemset lattice contains all the possible itemsets for a transaction database.

1 - itemsets A, B, C.



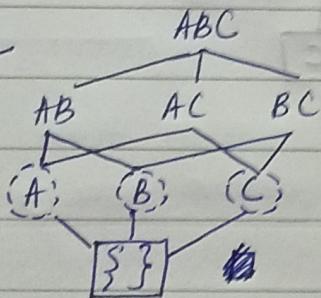
$AB \sqcup AC$   
 $n=1$

$AB \sqcup BC 1$

algorithm 3.

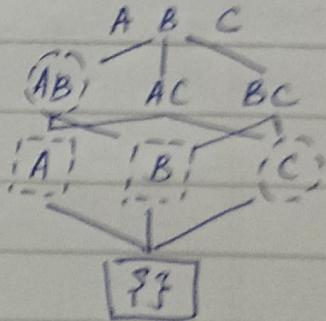
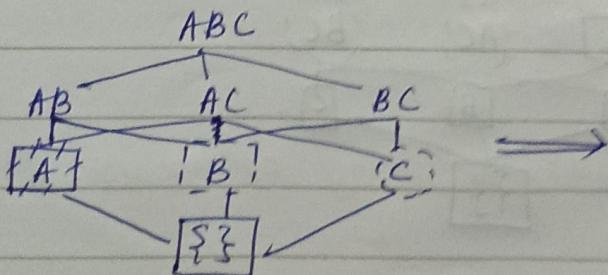
By Once a dashed <sup>itemset</sup> has been counted through all the transactions make it solid and stop counting it.

Step 1:



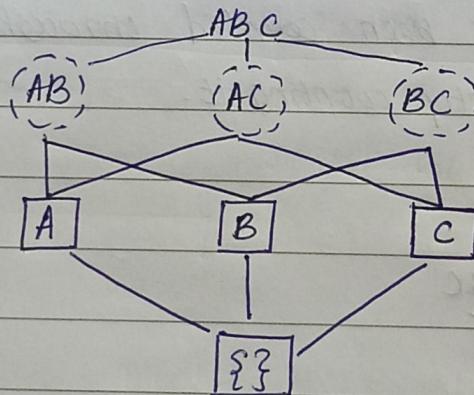
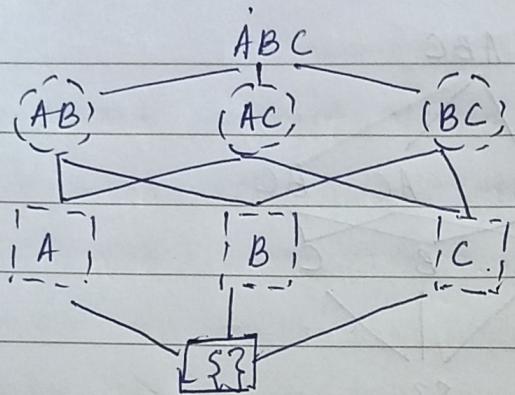
$M=2 \quad A \quad B \quad C$

T1 & T2      2      1      0

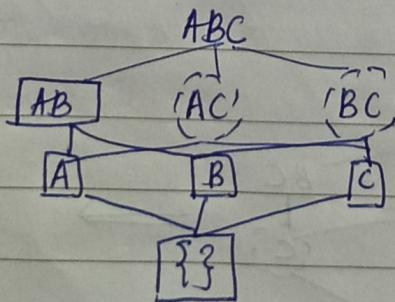


A      B      C  
M = 4      2      2      1

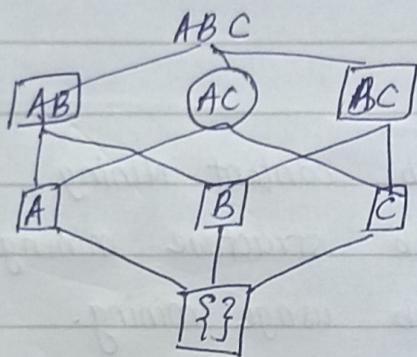
73 & 74



M = 6 (3M)



After 4M transactions,



## Web Mining

Mining of data related to the world wide web.

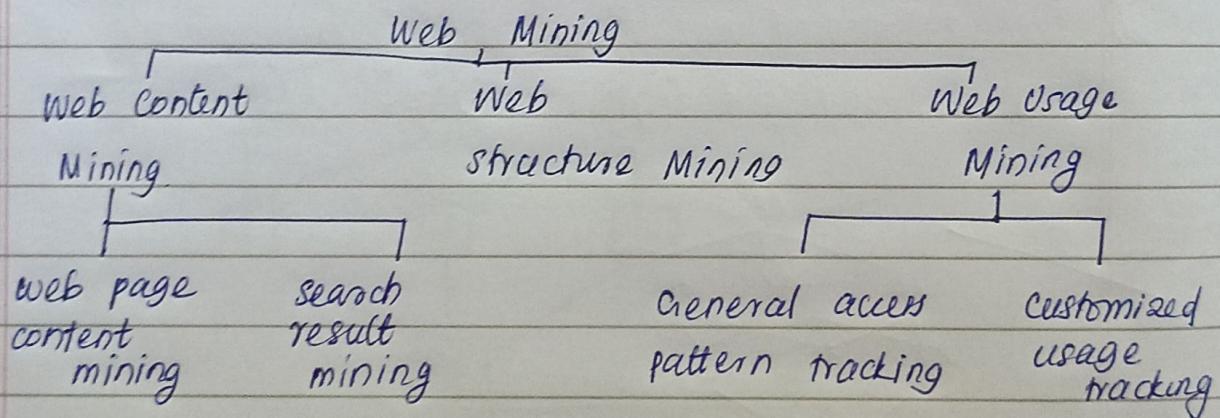
Used to discover patterns, structures and knowledge from the web.

3 areas : Web content Mining

Web structure mining

Web usage mining.

## Web Mining Taxonomy



- Q. T1 m, o, n, k, e, y
- T2 d, o, n, k, e, y
- T3 c, o, o, k, i, e
- T4 m, u, c, k, b, y

Find the frequent-path using FP-growth.

- Ans.
- T1  $\{ \overset{\checkmark}{e}, \overset{\checkmark}{k}, \overset{\checkmark}{m}, \overset{\checkmark}{n}, \overset{\checkmark}{o}, \overset{\checkmark}{y} \}$
- T2  $\{ d, e, k, n, o, y \}$
- T3  $\{ c, e, \cancel{k}, \cancel{i}, k, o, o \}$
- T4  $\{ \cancel{m}, \cancel{n}, c, k, m, u, y \}$

L1	Itemset	Sup-count	arrange in decreasing order of support count	Itemset	SC
<u>e</u>	2			<u>k</u>	4
<u>d</u>	1			<u>e</u>	3
<u>c</u>	3			<u>o</u>	3
<u>i</u>	1			<u>y</u>	3
<u>k</u>	4			<u>e</u>	2
<u>m</u>	2			<u>m</u>	2
<u>n</u>	2			<u>n</u>	2
<u>o</u>	3				
<u>u</u>	1				
<u>y</u>	3				

Transactions rearrange :

- T1  $\{ k, e, \overset{o}{m}, \overset{o}{y}, m, n \}$
- T2  $\{ k, d, e, n, o, y \}$   $\{ k, e, o, y, n, \cancel{d} \}$
- T3  $\{ k, c, e, i, o, o \}$   $\{ k, e, o, o, c, i \}$
- T4  $\{ k, c, m, u, y \}$   $\{ k, y, c, m, u \}$