

## UNIT - III

### Association Rules

Association rule learning is a type of unsupervised learning technique that checks for the dependency of one data item on another data item and maps accordingly so that it can be more profitable. It tries to find some interesting relations or association among the variable of dataset.

It is based on different rules to discover the interesting relations between variables in the database.

- Market basket analysis is a technique used by the various big retailer to discover the associations between items. We can understand it by taking an example of a supermarket as in a supermarket all products that are purchased together are put together.

For example if a customer buys bread, he most likely can also buy butter, eggs, or milk so these products are stored with in a shelf or mostly near by.

Customer 1	Customer 2	Customer 3	Customer 4
milk	milk	milk	Sugar
bread	bread	bread	eggs
butter	butter	butter	
	Cereals		

For example : market Based Data

Transaction ID	Items
1	{Milk, Bread, Rice, Book}
2	{Bread, Jam, Book, Pen}
3	{Jam, Milk, bread, Rice, Eggs}
4	{Rice, Eggs, Pen, Book}
5	{Eggs, Pen, Milk, Bread, Jam}
6	{Eggs, Rice, Bread, Jam}

Let us consider one transaction like

$$\{Milk, Bread, Rice, Book\} \quad \{Milk\} \rightarrow \{Bread\}$$
$$\{Bread, Jam, Book, Pen\} \quad \{Book\} \rightarrow \{Pen\} \quad \{Bread\} \rightarrow \{Jam\}$$

Some similar associations

$$\{Dishwash Liquid\} \rightarrow \{Scrubber\}$$
$$\{Laptop\} \rightarrow \{Mouse\}$$
$$\{Floor sticks\} \rightarrow \{Floor cleaner\}$$

Itemset : {Milk, Bread, Jam, Rice, Eggs, Book, Pen}

TID	Items						
	Milk	Bread	Jam	Rice	Eggs	Books	Pen
1	1	1			1	1	
2		1	1				1
3	1	1	1	1	1	1	
4	*	*		*	1	1	1
5	1	1	1		1		
6		1	1	1	1	1	

Frequent itemsets :-

Two itemsets : {Milk, Bread}, {Bread, Jam}, {Rice, Eggs}, {Books, Pen}

Three itemsets : {Milk, Bread, Jam}, {Rice, Eggs, Bread}, {Book, Pen, Egg}

Four itemsets : {Milk, Bread, Rice, Eggs}

Support :- It is a measure of how frequently a set of items occur in total no. of transactions.

$$\{Milk, Bread\} \rightarrow \{x, y\} \quad (x: Milk) (y: Bread)$$

Therefore The frequency of occurrence of x and y together in total no. of transactions is support -

$$\{Milk, Bread, Jam\} \rightarrow \{x, y\} \quad (x: Milk) (y: Bread, Jam)$$

Here the frequency of occurrence of {Bread, Jam} with {Milk} in whole transaction is support.

$$\text{Support}(S) = \frac{\sigma(x \cup y)}{N}$$

Confidence: It is a measure of how often items in Y appears in transactions that contain X.

$$\{milk, Bread, Jam\} \rightarrow \{x, y\} [x: milk] \{y: Bread, Jam\}$$

Here base the frequency of occurrence of x and y in all the transactions where x exists.

$$\text{Confidence}(C) = \frac{\sigma(x \cup y)}{\sigma x}$$

Association Rule Mining: Given a set of transactions T, the goal of association rule mining is to find all rules having

$\text{support} \geq \text{minsup}$  Threshold  
 $\text{confidence} \geq \text{minconf}$  Threshold

Eg: Suppose  $\text{minsup} = 0.3$   
 $\text{minconf} = 0.6$

$$\text{Consider: } \{Rice, Eggs\} \rightarrow \{x, y\}$$

$$\text{Then } \text{Support}(S) = \frac{\sigma(x, y)}{N}$$

$$\text{Support}(S) = \frac{1+1+1}{6} = \frac{3}{6} = 0.5$$

$$\text{and Confidence}(C) = \frac{\sigma(x \cup y)}{\sigma x}$$

$$= \frac{3}{4} = 0.75$$

Association Rule Mining.

here support : 0.5  $\geq$  minsup(0.3)

Confidence: 0.75  $\geq$  minconf(0.6)

Therefore we can mine {Rice, Eggs} as a rule.

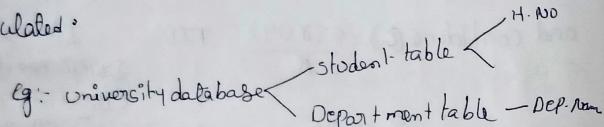
TID	Items
1	{Milk, Bread, Rice, Book}
2	{Bread, Jam, Book, Pen}
3	{Jam, Milk, Bread, Rice}
4	{Rice, Eggs, Pen, Book}
5	{Eggs, Pen, Milk, Bread, Jam}
6	{Eggs, Rice, Bread, Jam}

## Frequent Itemset Generation:-

Association mining searches for frequent items in the data-set.

In frequent mining usually the interesting associations and correlations between itemsets in transactional and relational databases are found.

**Dataset** :- A collection of related sets of information that is collected of separate elements but can be manipulated.



Frequent mining shows which items appear together in a transaction or relation.

### Need of Association Mining :-

- Frequent mining is generation of association rules from a transactional dataset.
- If there are 2 items X and Y purchased frequently then it's good to put item together in stores or provide some discount offer on one item on purchase of another item.
- This can really increase the sales.

For example it is likely to find that if a customer buys milk and bread he also buys butter.

$$\text{Support} := \frac{G(A \cup B)}{G}$$

$$\text{Confidence} := \frac{G(A \cup B)}{G(A)}$$

Support-Count(A) :- Number of transactions in which A appears. If A is ~~only~~ then it is the no. of transactions in which A and B are present.

Maximal Itemset :- An itemset is maximal frequent if none of its superset are frequent.

Closed Itemset :- An item set is closed if none of its immediate superset have same support count same as Itemset.

K-Itemset :- Item set which contains k items is a K-itemset

so it can be said that an itemset is frequent if the corresponding support count is greater than minimum support count.

Eg: Finding frequent items.

Transaction ID	Items
1	(A, C, D)
2	(B, C, D)
3	(A, B, C, D)
4	(B, D)
5	(A, B, C, D)

The APRIORI Principle :-

Apriori algorithm is a sequence of steps to be followed to find the most frequent itemset in the given database. Apriori algorithm is given by R. Agrawal and R. Srikant in 1994 for finding frequent itemsets in a dataset for boolean association rule. Name of the algorithm is Apriori because it uses prior knowledge of frequent itemset properties.

We apply an iterative approach or level-wise. Search where k-frequent itemsets are used to find k+1 itemsets.

To improve the efficiency of level wise generation of frequent itemsets, an important property which helps by is called Apriori Property which helps by reducing the search space.

Eg:

Transaction ID	List of items
1	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>
2	I <sub>2</sub> , I <sub>4</sub>
3	I <sub>2</sub> , I <sub>3</sub>
4	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>
5	I <sub>1</sub> , I <sub>3</sub>
6	I <sub>2</sub> , I <sub>3</sub>
7	I <sub>1</sub> , I <sub>3</sub> , I <sub>5</sub>
8	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>
9	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>

minimum support count is 2.  
minimum confidence is 67%.

Step 1: K = 1

(i) Create a table containing support count of each item present in dataset called C, (candidate set)

Itemset	Support-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

(ii) Compare Candidate set items support count with minimum support Count (here min-support = 2). If support-count of Candidate set items is less than min-support then remove those items. This gives us itemset  $L_1$ .

itemset	sup-count
$I_1$	6
$I_2$	7
$I_3$	6
$I_4$	2
$I_5$	2

$L_1$

Step 2:  $k=2$

Generate Candidate set  $C_2$  using  $L_1$

itemset	sup-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

$C_2$

(ii) Compare Candidate ( $C_2$ ) support count with minimum support count ( $= 2$ ). If support count of Candidate set item is less than min-support then remove Those items) this gives us itemset  $L_2$

itemset	sup-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

~~$I_{3,4}$~~

$L_2$

Step-3:- Generate Candidate set  $C_3$  using  $L_2$

itemset	sup-count
$I_1, I_2, I_3$	2
$I_1, I_2, I_5$	2
$I_1, I_3, I_5$	1
$I_1, I_2, I_4$	1

Compare Candidate set items with min-sup-count = 2, if support count is less than min-supcount then remove Those items. This gives  $L_3$

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

In L2 These sets are Pragent

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

In L2 These sets are Pragent

Confidence :-

A confidence of 60% means That 60% of the customers, who purchased m

Confidence :-

A confidence of 60% means That 60% of the customers, who purchased milk and bread also bought butter.

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

so here by taking an example of any frequent itemset, we will show the rule generation

Itemset  $\{I_1, I_2, I_3\}$  from L3

so rules can be

$$[I_1 \wedge I_2] \Rightarrow I_3 \text{ confidence} =$$

$$\frac{\text{sup}[I_1 \wedge I_2 \wedge I_3]}{\text{sup}[I_1 \wedge I_2]} = \frac{2}{4} * \frac{100}{2} = 0.50 * 50\%$$

$$[I_1 \wedge I_3] \Rightarrow [I_2]$$

$$\frac{\text{sup}[I_1 \wedge I_2 \wedge I_3]}{\text{sup}[I_1 \wedge I_3]} = \frac{2}{4} * 100 = 50\%$$

$$[I_2 \wedge I_3] \Rightarrow [I_1]$$

$$\frac{\text{sup}[I_1 \wedge I_2 \wedge I_3]}{\text{sup}[I_2 \wedge I_3]} = 2/4 * 100 = 50\%$$

$$[I_1] \Rightarrow [I_2 \wedge I_3]$$

$$\frac{\text{sup}[I_1 \wedge I_2 \wedge I_3]}{\text{sup}[I_1]} = \frac{2}{6} * 100 = 33\%$$

$$[I_2] \Rightarrow [I_1 \wedge I_3]$$

$$\frac{\text{sup}[I_1 \wedge I_2 \wedge I_3]}{\text{sup}[I_2]} = \frac{2}{7} * 100 = 28\%$$

$$[I_3] \Rightarrow [I_1 \wedge I_2]$$

$$\frac{\text{sup}[I_1 \wedge I_2 \wedge I_3]}{\text{sup}[I_3]} = \frac{2}{6} * 100 = 33\%$$

so if minimum confidence is 50%. Then first 3 rules can be considered as strong association rules.

Another Example :-

minimum support count = 2

TID List of items

T10  $I_1, I_3, I_4$

T20  $I_2, I_3, I_5$

T30  $I_1, I_2, I_3, I_5$

T40  $I_2, I_5$

## The Partition Algorithms :-

This clustering method classifies the information into multiple groups based on the characteristics and similarity of the data. It's the data analysts to specify the number of clusters that has been generated for the clustering methods.

In the Partition method when Database (D) that contains multiple ( $N$ ) objects then the partitioning methods, constant user specified ( $K$ ) partitions of the data in which each partition represents a cluster and a particular region.

There are many algorithms that come under partitioning methods some of the popular ones are K-Mean, PAM (K-Medoids), CLARA algorithm (cluster large applications) etc.

## K-Mean clustering algorithm :-

The k-mean algorithm takes the input parameters  $K$  from the user and partitions the data set containing  $N$  objects into  $K$  clusters so that the resulting similarity among the data objects inside the group (intracluster) is high but the similarity of data objects with the data objects from the outside the cluster is low (intercluster).

The similarity of the cluster is determined with respect to the mean value of the cluster

In K-mean clustering algorithm we should follow some steps

Step 1: Take the Mean value.

Step 2: Find the nearest no. of mean and put in Cluster.

Step 3: Repeat ① & ② until we get the same mean

### Example :-

$$\text{Data Pth, } K = \{2, 4, 6, 9, 18, 16, 20, 24, 26\}$$

No. of clusters = 2.

Take two Random numbers as clusters = {4, 12} Centroids

$$K_1 = \{2, 4, 6\} \quad K_2 = \{9, 12, 16, 20, 24, 26\}$$

$$\text{Mean value for } K_1 = \frac{2+4+6}{3} = 4 \quad \text{Mean value for } K_2 = \frac{9+12+16+20+24+26}{6}$$

$$K_2 = 18(17.8)$$

$$K_1 = \{2, 4, 6, 9\}$$

$$K_2 = \{12, 16, 20, 24, 26\}$$

$$\text{Mean } K_1 = \frac{21}{4} = 5.25 \\ = 5$$

$$\text{Mean } K_2 = 19.6 \\ = 20$$

$$K_1 = \{2, 4, 6, 9, 12\}$$

$$K_2 = \{16, 20, 24, 26\}$$

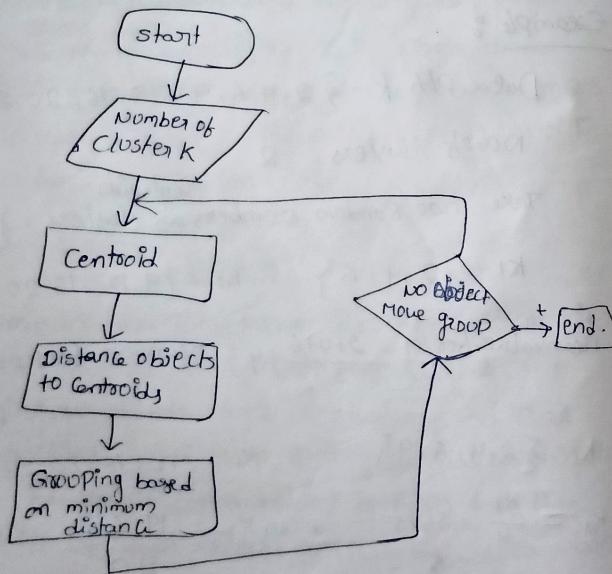
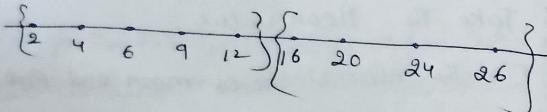
$$\text{Mean value } K_1 = 6.6$$

$$K_2 = 21.5 = 22$$

$$K_1 = \{2, 4, 6, 9, 12\} \quad K_2 = \{16, 20, 24, 26\}$$

mean value  $K_1 = 6.6 = 7$

mean value  $K_2 = 24.5 = 24$



Divide the given sample data in two (2) clusters using k-means algorithm [Euclidean distance]

	Height	Weight
1	185	72
2	170	56
3	168	60
4	179	68
5	182	72
6	188	77
7	180	71
8	180	70
9	183	84
10	180	88
11	180	67
12	177	76

$$\sqrt{(x_H - H_i)^2 + (x_W - W_i)^2}$$

↓

Observed value      Centroid value      Observed value

## FP-Growth algorithms:-

This Algorithm is an improvement to the Apriori Algorithm. A Frequent Pattern (FP) is generated without the need for Candidate generation. FP Growth algorithm represents the database in the form of a tree called a frequent pattern tree or FP-tree.

This tree structure will maintain the association between the itemsets. The database is fragmented using one frequent item. This Fragmented Part is called "pattern fragment". The itemset of these fragmented patterns are analyzed. This will reduce the search for frequent itemsets comparatively.

## FP-Tree:-

Frequent pattern tree is a tree-like structure that is made with the initial itemsets of the database. The purpose of the FP tree is to mine the most frequent pattern. Each node of the FP-tree represents an item of the itemset.

The root node represents null while the lower nodes represent the itemsets. The association of the nodes with the lower nodes that is the itemsets with the itemsets are maintained while forming the tree.

## Frequent Pattern Algorithm steps :-

The frequent pattern growth method lets find the frequent pattern with out candidate generation.

1. The first step is to scan the database to find the occurrences of the itemsets in the database. This step is the same as the first step of Apriori. The count of one itemset in the database is called support count or frequency of 1-itemset.
2. The second step is to construct the FPtree. For this, create the root of the tree. The root is represented by null.
3. The next step is to scan the database again and examine the transactions. Examine the first transaction and find out the itemset in it. The item set with the max count is taken at the top. The next itemset with lower count and so on. It means that the branch of the tree is constructed with transaction itemsets in descending order of count.

4) The next transaction in the database is examined. The itemsets are ordered in descending order of count. If any itemset of this transaction is already present in another branch then this transaction branch

would share a common prefix to the root. This means that the common itemset is linked to the new node of another itemset in this transaction.

5) Also, the count of the itemset incremented as it occurs in the transactions. Both the common node and new node count is increased by 1 as they are created and linked according to transactions.

6) The next step is to mine the created FPtree.

For this, the lowest node is examined first along with the links of the lowest nodes.

The lowest node represents the frequent pattern length 1. From this traverse the path in the

FP tree. This path or paths called a conditional pattern base. Conditional pattern base is a sub-database consisting of prefix paths in the FP tree occurring with the lowest node.

7. Construct a conditional FPtree, which is formed by a count of itemsets in the path. The itemsets meeting the threshold support are considered in the conditional FP tree.
8. Frequent patterns are generated from the conditional FP tree.

example:

FP-Tree for the following transaction Data set [minimum support = 30%]

minimum no. of transactions =  $0.3 * 8 =$

$2.4 \approx 3$

Trans ID	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

Lower priority number means higher priority.

order the items according to priority



## Compact Representation of Frequent Item Set

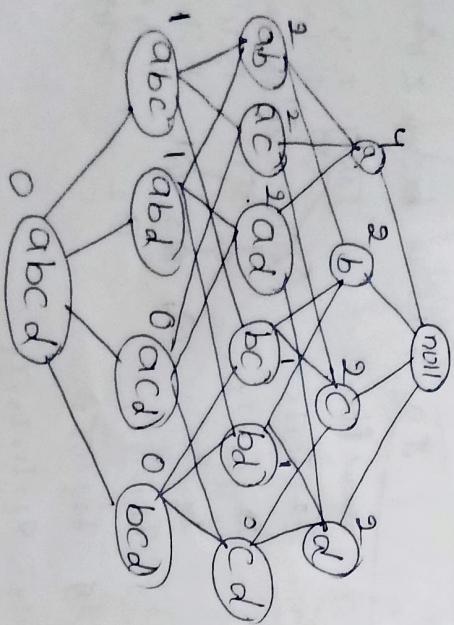
The no of frequent itemsets produced from a transaction dataset can be very large. It is useful to identify a small representative set of itemsets from which all other frequent itemsets can be derived.

In this, we have two types.

### i) maximal frequent itemset

A maximal frequent itemset is defined as a frequent itemset for which none of its immediate supersets are frequent.

Ex:-



The support counts are shown on the top left of each node. Assume support count threshold = 50%. It is each itemset must occur in 2 or more transactions. Based on that threshold, the frequent itemsets are a, b, c, d, ab, ac, and ad.

Out of the 7 frequent itemsets 3 are identified as maximal frequent itemsets.

ab: immediate supersets abc and abd are infrequent

ac: immediate superset abc and acd are infrequent -

ad: immediate supersets abd and acd are infrequent.

The remaining 4 frequent nodes (a, b, c, d) cannot be maximal frequent because they all have at least 1 immediate superset that is frequent.

## Advantages:

Maximal frequent itemsets provide a compact representation of all the frequent itemsets for a particular dataset.

## Disadvantages:

The support count of maximal frequent itemsets does not provide any information about the support count of their subsets. This means that an additional traversal of data is needed to determine support count for non-maximal frequent itemsets.

## Frequent itemset:

A frequent itemset is an itemset whose support is greater than some user-specified minimum support.

## Closed Frequent itemset:

An itemset is closed if none of its immediate supersets has the same support as that of the itemset.

## Maximal frequent itemset:

An itemset is maximal frequent if none of its immediate superset is frequent.

→ Need of closed and maximal itemsets:-

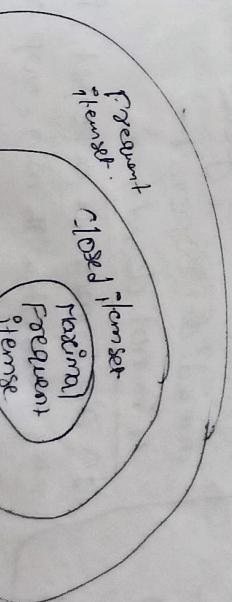
→ These closed frequent itemsets and maximal frequent itemsets are useful when huge amount of data is used in association rule mining.

→ If the length of frequent itemset is 'k' then by downward closure property, all of its  $2^k$  subsets are also frequent.

→ When the computation is very expensive there is no interest to find additional subsets. This can be avoided by frequent itemset with maximum length.

→ One disadvantage with maximal frequent itemsets is that even all its subsets are frequent, we do not know their supports. For mining rules, support information is very important.

→ These type of closed frequent itemset is preferred.



g)

Item	list of items
T <sub>1</sub>	A, B, C, D
T <sub>2</sub>	A, B, C, D
T <sub>3</sub>	A, B, C
T <sub>4</sub>	B, C, D
T <sub>5</sub>	C, D

min-sup-count = 3

support set

B is not closed.

In B's immediate superset, itemset are present

min min-sup-count 3.

B is not maximal.

C(5) → A&(3), BC(4), CD(2)

C(count) is greater than its immediate supersets.

C is closed.

In C's immediate superset, itemsets are present with minimum support count 3.

C is not maximal.

Item	count
A	3
B	4
C	5
D	4

All items that is AB,BC,CD

are frequent because they

Support count is greater than

as equal to minimum sup-count.

D(4) → AD(2), BD(3), CD(4)

D(count) is not greater than its immediate superset

D is not closed.

In D's immediate superset, itemset are present with

min sup.count 3.

D is not maximal.

A is not maximal.

B(4) → AB(3), BC(4), BD(3).

B(count) is not greater than its immediate superset

Item	Count
A, B, C	3
A, B, D	2
A, C, D	2
B, C, D	3

$ABC \rightarrow ABC(3), ABD(2)$

$BCD(3) \rightarrow ABCD(2)$

$BCD$  is closed.

$BCD$  is maximal.

$AB$  is immediate superset, Itemset are present with min. support count 3

$AB$  is not maximal.

$AC(3) \rightarrow ABC(3), ACD(2)$

$AC$  is not closed

$AC$  is not maximal.

$AD(2) \rightarrow ABD(2), ACD(2)$

$AD$  not frequent

$BC(4) \rightarrow ABC(3), BCD(3)$

$BC$  is closed

$BC$  is not maximal

$BD(3) \rightarrow ABD(2), BCD(3)$

$BD$  is not closed

$BD$  is not maximal

$CD(4) \rightarrow ACD(2), BCD(3)$

$CD$  is closed

$CD$  is not maximal

$ABC(3) \rightarrow ABCD(2)$

$ABC$  is closed

$ABC$  is maximal.

$ABD(2) \rightarrow ABD(2) \rightarrow$  not frequent

$ABD(2) \rightarrow ABD(2) \rightarrow$  not maximal.

Item	Count
A, B, C, D	2

1-item set			
Item	Freq	Closed	Maximal
A(3)	✓	✗	✗
B(4)	✗	✗	✗
C(5)	✗	✗	✗
D(4)	✗	✗	✗

2-item set			
Item	Freq	Closed	Max
ABC(3)	✗	✗	✗
ABD(2)	✗	✗	✗
ACD(2)	✗	✗	✗
BCD(3)	✗	✗	✗
CD(4)	✗	✗	✗

3-item set			
Item	Freq	Closed	Maximal
ABC(3)	✗	✗	✗
ABD(2)	✗	✗	✗
ACD(2)	✗	✗	✗
BCD(3)	✗	✗	✗

4-item set			
Item	Freq	Closed	Maximal
ABCD(2)	✗	✗	✗

## TID List of items

T <sub>1</sub>	B, J, P	Min-SUP count = 40%
T <sub>2</sub>	B, P	find frequent, closed, maximal
T <sub>3</sub>	B, M, P	itemsets
T <sub>4</sub>	E, B	
T <sub>5</sub>	E, M	

Class

Input

gt is the

Predefined

gt is the

maps each

classification

## Purpose of classification

Descriptive

A classification

approach to

data set. gt

In classification

a training da

algorithm deri

model can be a d

In test data set

model. Then if