

Association Rule Mining

→ It discovers the probability of occurrence of items in a collection. helps in discovering some interesting relationships in large data sets.

→ A data set contains data objects and each data object contains a set of attributes. An attribute is called as dimension or feature or variable which represents the characteristic feature of the data object.

eg: Height, qualification, colour etc.

Association Rule Mining

It finds the interesting associations and relationships among large sets of data items. This rules shows how frequently a itemset occurs in a transaction.

for eg: Market Basket Data

Trans ID	Item
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 →	{ Egg, Rice, Bread, Jam }

Let us consider one transaction like.

{Milk, Bread, Rice, Book}

{Milk} \rightarrow {Bread}

{Bread, Jam, Book, Pen}

{Book} \rightarrow {Pen}

{Bread} \rightarrow {Jam}

Interesting patterns

Some Similar associations

{Milk with liquid} \rightarrow {scrubber}

{Laptop} \rightarrow {Mouse}

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

Frequent Item sets:

Two Itemsets: {M, B}, {B, Jam}, {Rice, Eggs}, {Book, Pen}

Three Itemset: {M, B, Jam}, {Rice, Eggs, Bread}, {Book, Pen, Eggs}

Four Itemset: {Milk, Bread, Rice, Eggs} etc.

sq, can, rule mining \rightarrow 5 steps.

Support: It is a measure of how frequently a set of items occur in total number of transactions.

eg: {Milk, bread} \rightarrow {X, Y} . {X \rightarrow Milk} {Y \rightarrow Bread}

Therefore the frequency of occurrence of X and Y together in total no of transactions is support.

eg {Milk, Bread, Jam} \rightarrow {X, Y} . {X \rightarrow Milk}, {Y \rightarrow Bread, Jam}

\Rightarrow Here the frequency of occurrence of {Bread, Jam} with {Milk} in whole transaction is support. Support is $= \frac{\sigma(X \cup Y)}{N} =$

Confidence is It is a measure of how often item y appears in transactions that contain x .

$\{ \text{Milk, Bread (T=1)} \} \rightarrow (x, y)$

$(x \rightarrow \text{Milk}) \mid (y : \text{Bread (T=1)})$

Therefore 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 ,
goal of association rule mining is to find all rules having

$$\text{support} \geq \text{min sup threshold}$$
$$\text{confidence} \geq \text{min conf threshold}.$$

Eg:- Suppose, $\text{min sup} = 0.3$
 $\text{min conf} = 0.6$

Consider ^{2-Items} $\{ \text{Rice, Eggs} \} \rightarrow \{x, y\}$

Then

Support (S) =	$\frac{\sigma(xy)}{n}$
---------------	------------------------

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

and

$$\text{confidence}(C) = \frac{\sigma(xy)}{\sigma x}$$
$$= \frac{3}{4} = 0.75$$

Association Rule Mining [Threshold Value]

here $\text{support} = 0.5 \geq \text{min sup} (0.3)$

$\text{confidence} = 0.75 \geq \text{min conf} (0.6)$

Therefore, we can mine.

$\{ \text{Rice, Eggs} \}$ as ~~a rule~~ a rule.

(2)

eg 2:- suppose

$$\text{minsup} = 0.3$$

$$\text{minconf} = 0.6$$

Consider, {Milk, Bread, Jam} \rightarrow {X, Y}

$$X = \{\text{Milk}\} \quad Y = \{\text{Bread, Jam}\}$$

then $\text{support}(s) = \frac{\sigma(X \cup Y)}{N}$

$$= \frac{2}{6} = 0.333$$

and $\text{confidence}(c) = \frac{\sigma(X \cup Y)}{\sigma X}$

$$= \frac{2}{3} = 0.667$$

Association Rule Mining

here, $\text{support} : 0.333 \geq \text{minsup}(0.3)$

$\text{confiden} : 0.667 \geq \text{minconf}(0.6)$

Therefore we can mine {Milk, Bread, Jam} as a Rule.

" Apriori Algorithm :

Transaction list :-

1	Milk	Egg	Bread	Butter
2	Milk	Butter	Egg	Ketchup
3	Bread	Butter	Ketchup	
4	Milk	Bread	Butter	
5	Bread	Butter	Cookies	
6	Milk	Bread	Butter	Cookies
7	Milk	Cookies		
8	Milk	Bread	Butter	
9	Bread	Butter	Egg	Cookies
10	Milk	Butter	Bread	
11	Milk	Bread	Butter	
12	Milk	Bread	Cookies	Ketchup

- ⇒ Here are 12 dozen sale transactions
- ⇒ The objective is to use this transaction data to find affinities between products, that is, which products sell together often.
- ⇒ The support level will be set at 30% percent, the confidence level will be set at 50% percent.
- ⇒

Association Rule Mining

②

1-Item set	frequency
Milk	9
Bread	10
Butter	10
Egg	3
Ketchup	5
Cookie	5

requent-1-Itemset	frequency
Milk	9
Bread	10
Butter	10
Cookie	5

$\frac{4}{12} = 33\%$. Here egg, ketchup is not satisfied here min-support count
[4 above transaction's only
work here]

2-Item set	frequency
Milk, Bread	7
Milk, Butter	7
Milk, Cookie	3
Bread, Butter	9
Butter, Cookie	3
Bread, Cookie	4

requent 2-Itemset	frequency
Milk, Bread	7
Milk, Butter	7
Bread, Butter	9
Bread, Cookie	4

It should appear more than 4. Then then only it will
satisfying the min-support count.
here {Milk, Cookie, Butter, Cookie} → 3 [min-supp- 30%]

{ Milk, Bread, Butter, Cookie }

3-Item set	frequency
Milk, Bread, Butter	6
Milk, Bread, Cookie	1
Bread, Butter, Cookie	3
Milk, Butter, Cookie	2

one only.

requent 3-Itemset	frequency
Milk, Bread, Butter	6

Association Rule Mining - Subset Creation

→ frequent 1-Itemset set = $I \Rightarrow \{Milk, Bread, Butter\}$

→ Non-empty subset are

- $\{Milk\}, \{Bread\}, \{Butter\}, \{Milk, Bread\}, \{Milk, Butter\}, \{Bread, Butter\}$

→ How to form Association Rule - b.

→ For every non-empty subset S of I , the association

rule is,

$$S \rightarrow (I - S)$$

* If $\text{support}(I) / \text{support}(S) \geq \text{min-confidence}$.

Non-empty subset are:-

- $\{Milk\}, \{Bread\}, \{Butter\}, \{Milk, Bread\}, \{Milk, Butter\}, \{Bread, Butter\}$

min-support = 30% and confidence = 60%.

Rule 1:

$\{Milk\} \rightarrow$

$\{Bread, Butter\}$

$\{S = 50\%, C = 66.67\%$

- support = $6/12 = 50\%$

- confidence = $\frac{\text{support}\{Milk, Bread, Butter\}}{\text{support}\{Milk\}}$

$$= \frac{6/12}{9/12} = \frac{6}{9} = 66.67 \geq 60\%$$

- valid.

Rule 2:

$\{Bread\} \rightarrow$

$\{Milk, Butter\}$

$\{S = 50\%, C = 60\%$

- support = $6/12 = 50\%$

- confidence = $\frac{\text{support}\{Milk, Bread, Butter\}}{\text{support}\{Bread\}}$

$$= \frac{6}{10} = 60\% \geq 60\%$$

→ valid

(2)

Rule 3: $\{ \text{Butter} \} \rightarrow \{ \text{Milk, Bread} \}$ $\{ s=10\%, c=60\% \}$

- support = $6/12 = 50\%$

- confidence = $\text{support}(\text{Milk, Bread, Butter}) / \text{support}(\text{Butter})$

- valid $= \frac{6}{10} = 60\% \geq 60$

Rule 4: $\{ \text{Milk, Bread} \} \rightarrow \{ \text{Butter} \}$ $\{ s=10\%, c=85.7\% \}$

- support = $6/12 = 50\%$

- confidence = $\text{support}(\text{Milk, Bread, Butter}) / \text{support}(\text{Milk, Bread})$

- valid $= \frac{6}{7} = 85.7\% \geq 60\%$

Rule 5: $\{ \text{Milk, Butter} \} \rightarrow \{ \text{Bread} \}$ $\{ s=50\%, c=85.7\% \}$

- support = $6/12 = 50\%$

- confidence = $\text{supp}(\text{Milk, Bread, Butter}) / \text{support}(\text{Milk, Butter})$

$= \frac{6}{7} = 85.7\% \geq 60\%$

Rule 6: $\{ \text{Bread, Butter} \} \rightarrow \{ \text{Milk} \}$ $\{ s=50\%, c=66.6\% \}$

- support = $6/12 = 50\%$

- confidence = $\text{support}(\text{Milk, Bread, Butter}) / \text{support}(\text{Bread, Butter})$

$= \frac{6}{9} = 66.6\% \geq 60$

⇒ Frequent pattern (FP) Growth Algorithm

set of transactions

Transaction ID	Item
T ₁	{E, K, M, N, O, Y}
T ₂	{D, E, K, N, O, Y}
T ₃	{A, E, K, M}
T ₄	{C, K, M, O, Y}
T ₅	{C, E, I, K, O, Y}

few items are bought together

Item	frequency
A	1
C	2
D	1
E	4
I	1
K	5
M	3
N	2
O	3
Y	3

- ⇒ The above-given data is a hypothetical dataset of transactions with each letter representing an item.
- ⇒ Let the minimum support be 3.
- ⇒ A frequent pattern set (L)₀ is built which will contain all the elements whose frequency is greater than or equal to the minimum support.
- As minimum support is 3.
- These elements are stored in descending order of their respective frequencies.
- After insertion of the relevant items, the set L looks like this :-
- $$L = \{K:5, E:4, M:3, O:3, Y:3\}$$
- Decreasing order

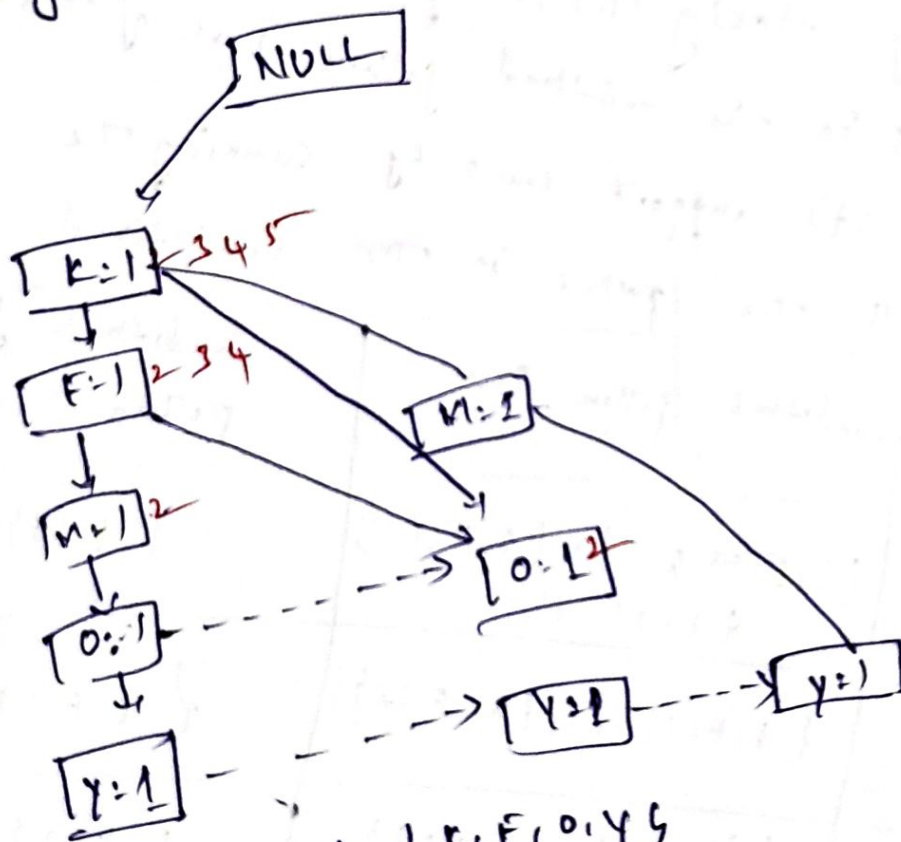
Now, for each transaction, the respective Order-Item set is built.

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

Transaction ID	Item	Ordered-Item set
T ₁	{E, K, M, N, O, Y}	{K, E, M, O, Y}
T ₂	{D, E, K, N, O, Y}	{K, E, O, Y}
T ₃	{A, E, K, M}	{K, E, M}
T ₄	{C, K, M, O, Y}	{K, M, Y}
T ₅	{C, E, I, K, O, Y}	{K, E, O}

⇒ Now, all the Ordered-Itemset are Inserted into a Trie Data Structure.

a) Inserting the set {K, E, M, O, Y}



b) Inserting the set {K, E, O, Y}

c) {K, E, M}

d) {K, E, Y}

e) {K, E, O}

Now for each item, the Conditional pattern base is computed which is path labels of all the paths which lead to node of the given item in the frequent pattern tree.

Item	Conditional pattern Base
Y	$\{\{K, E, M, O:1\}, \{K, F, O:1\}, \{K, M:1\}\}$
O	$\{\{K, E, M:1\}, \{K, F:2\}\}$
M	$\{K, E:2\}, \{K:1\}$
E	$\{K:4\}$
K	

Now for each item the Conditional frequent pattern Tree is built: It is done by taking the set of elements which is common in all the paths in the conditional pattern base of that item and calculating its support count by summing the support count of all the paths in the Conditional pattern base.

Item	Conditional pattern base	Conditional frequent pattern Tree
Y	$\{\{K, E, M, O:1\}, \{K, F, O:1\}, \{K, M:1\}\}$	$\{K:3\}$
O	$\{\{K, E, M:1\}, \{K, F:2\}\}$	$\{K, F:3\}$
M	$\{K, E:2\}, \{K:1\}$	$\{K:3\}$
E	$\{K:4\}$	$\{K:4\}$
K		

from the conditional frequent pattern tree, the frequent pattern entries are generated by pairing the items of the conditional frequent pattern tree set to the corresponding item as given in the below table.

Item	Frequent pattern generated
Y	$\langle K, Y:1 \rangle$
O	$\{ \langle K, O:3 \rangle, \langle E, O:3 \rangle, \langle E, K, O:3 \rangle \}$
M	$\{ \langle K, M:1 \rangle \}$
E	$\{ \langle E, K:2 \rangle \}$
K	

the association rules.