

## Association Rule Mining :-

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 associations among the variables of datasets. The association rule learning is one of the very important concept of machine learning, and it is employed in market based Analysis, Web usage mining, continuous production etc.

For example. If a customer buys bread he most likely can also buy butter, eggs, or milk. So these products are stored within a shelf or mostly nearby.

Association rule learning can be divided into three type of algorithms:-

- 1) Apriori
2. Eclat
- 3 F-P Growth Algorithm

Apriori Algorithm It uses association rules & it is designed to work on datasets or databases that contain transactions.

This algo. uses a BPS & Hash tree to calculate the itemset associations efficiently.

Frequent itemset — FI are those item whose support is greater than threshold value.

Eg:-  $A = \{1, 2, 3, 4, 5\}$   $B = \{2, 3, 7\}$

2, 3 are FI values.

Step-I Determine support of itemsets & Select min support & confidence.

Step-II Take all supports in transaction with higher support value than min or selected value.

Step-III Find all rules of these sets that have higher confidence value than threshold min confidence

Step-IV Sort the rules as decreasing order of lift.

## Apriori Algorithm:-

It uses frequent itemsets to generate association rules & it is designed to work on datasets or databases that contain transactions.

This algo. uses a BPS & Hash tree to calculate the itemset associations efficiently.

Frequent itemset :- FI are those items whose support is greater than threshold value.

Eg:-  $A = \{1, 2, 3, 4, 5\}$      $B = \{2, 3, 7\}$

2,3 are FI values.

Step-I Determine support of itemsets & Select min support & confidence.

Step-II Take all supports in transaction with higher support value than min or selected support value.

Step-III Find all rules of these sets that have higher confidence value than threshold or min confidence

Step-IV Sort the rules as decreasing order of lift.

# Apriori Algo.

②

Eg:-

Tid.	Itemsets
T <sub>1</sub>	A, B
T <sub>2</sub>	B, D
T <sub>3</sub>	B, C
T <sub>4</sub>	A, B, D
T <sub>5</sub>	A, C
T <sub>6</sub>	B, C
T <sub>7</sub>	A, C
T <sub>8</sub>	A, B, C, E
T <sub>9</sub>	A, B, C

Given S

min support = 2.

min confidence = 50%.

C → Customer  
FI → Frequent itemset

Step ① Calculating C<sub>1</sub> and F<sub>1</sub>

Itemset	Support count
A	6
B	7
C	6
D	2
E	1

→ F<sub>1</sub>

Itemset	SC
A	6
B	7
C	6
D	2

Step ② Calculating C<sub>2</sub> and F<sub>2</sub>.

Itemset	SC
(A, B)	4
(A, C)	4
(A, D)	1
(B, C)	4
(B, D)	2
(C, D)	0

→ F<sub>2</sub>

Itemset	SC
(A, B)	4
(B, C)	4
(A, C)	4
(B, D)	2

Step ③ Candidate Generation  $C_3 \& F_3$

is sc	
(A, B, C)	2
(B, C, D)	0
(A, C, D)	0
(A, B, D)	1

$F_3$

Itemset	Support
(A, B, C)	2

Step 4) Finding Association rules for subsets:-

→ We will calculate the confidence using

$$\frac{\text{Sup}(A \cap B)}{A} \rightarrow \frac{\text{Sup}(A \cap B \cap C)}{\text{Sup}(A \cap B)} \cdot \text{Sup} = \frac{A \rightarrow B}{A}$$

→ After calculating we will exclude rules that have less ~~than~~ confidence than min threshold (50%).

Rules	Support	Confidence.
$A \cap B \rightarrow C$	2	$\text{Sup}((A \cap B) \cap C) / \text{Sup}(A \cap B) = 2/4 = 50\%$
$B \cap C \rightarrow A$	2	$\text{Sup}((B \cap C) \cap A) / \text{Sup}(B \cap C) = 2/4 = 50\%$
$A \cap C \rightarrow B$	2	$\text{Sup}((A \cap C) \cap B) / \text{Sup}(A \cap C) = 2/4 = 50\%$
$C \rightarrow A \cap B$	2	$\text{Sup}((C \cap (A \cap B)) / \text{Sup}(C) = 2/5 = 40\%$
$A \rightarrow B \cap C$	2	$\text{Sup}((A \cap (B \cap C)) / \text{Sup}(A) = 2/6 = 33.33\%$
$B \rightarrow A \cap C$	2	$\text{Sup}((B \cap (A \cap C)) / \text{Sup}(B) = 2/7 = 28.57\%$

Frequent Pattern Growth Algo! → It is improved version of Apriori Algo. It represents the database in the form of tree structure. that is known as frequent pattern or tree.

Tid	Item
T <sub>1</sub>	E, K, M, N, O, Y
T <sub>2</sub>	D, E, K, N, O, Y
T <sub>3</sub>	A, E, K, M
T <sub>4</sub>	C, K, M, U, Y
T <sub>5</sub>	C, E, I, K, O, O

min Support = 3.

②

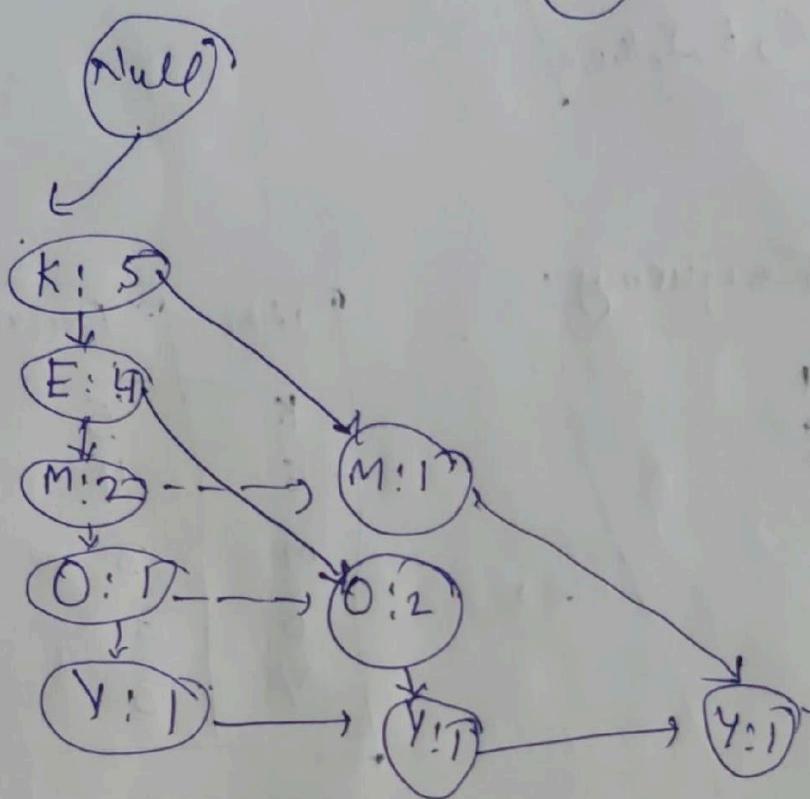
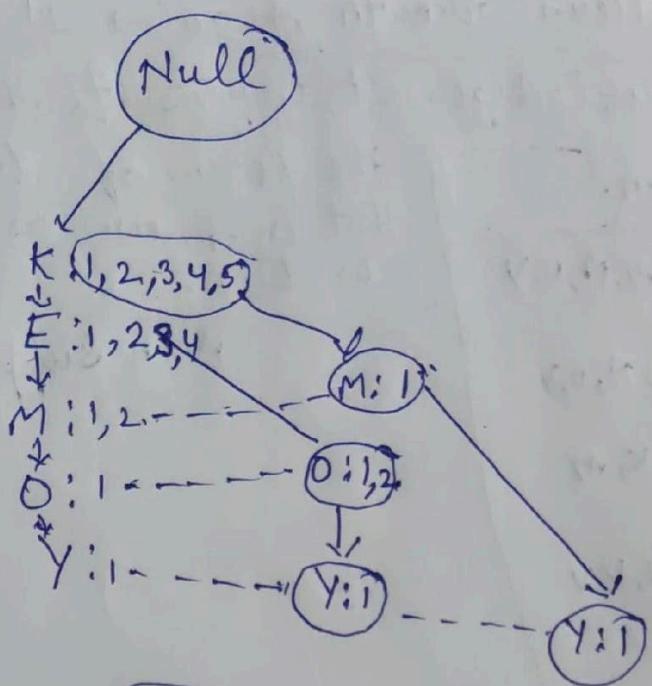
Item	Frequency
A	1
C	2
D	1
E	4
I	1
K	5
M	3
N	2
O	3
U	1
Y	3



Item	Frequency
K	5
E	4
M	3
O	3
Y	3

③

Tid	Items	Needed - Itemset
T <sub>1</sub>	E, K, M, N, O, Y	K, E, M, O, Y
T <sub>2</sub>	D, E, K, N, O, Y	K, E, O, Y
T <sub>3</sub>	A, E, K, M	K, E, M
T <sub>4</sub>	C, K, M, U, Y	K, M, Y
T <sub>5</sub>	C, E, I, K, O, O	K, E, O



Item.	Conditional Pattern Base	Conditional Freq.
Y.	{K, E, M, O: {1}}, {K, E, O: {1}}, {K, M, Y: {1}}	3K: 33
O	{K, E, M: {1}}, {K, E: {2}}	3K, E: 33
M.	{K, E: {2}}, {K: {1}}	3K: 33
E	{K: {4}}	3K: 43

Items

Frequent Pattern Generated

(4)

Y	$\{<K, Y : 3>\}$
O	$\{<K, O : 3>, <E, O : 3>, <K, E, O : 3>\}$
M	$\{<K, M : 3>\}$
E	$\{<K, E : 3>\}$
K	—

 $K \rightarrow Y, Y \rightarrow K$  $K \rightarrow O, E \rightarrow O, O \rightarrow K, O \rightarrow E, E \rightarrow K, K \rightarrow E$  $K \rightarrow M, M \nrightarrow K$  $K \rightarrow E, E \rightarrow K$ 

### \* Eclat Algorithm

Eclat Algorithm stands for Equivalence class clustering and bottom-up lattice Traversal.

It is one of popular methods of Association rule mining Algorithm.

Transaction id.	Bread	Butter	Milk	Coke	Jam
T <sub>1</sub>	1	1	0	0	1
T <sub>2</sub>	0	1	0	1	0
T <sub>3</sub>	0	1	1	0	0
T <sub>4</sub>	1	1	0	1	0
T <sub>5</sub>	1	0	1	0	0
T <sub>6</sub>	0	1	1	0	0
T <sub>7</sub>	1	0	1	0	0
T <sub>8</sub>	1	1	1	0	1
T <sub>9</sub>	1	1	1	0	1

$K=1$ , Minimum Support = 2.  $K=2$ , MS = 2.

Item	Tidset
Bread	$T_1, T_4, T_5, T_7, T_8, T_9$
Butter	$T_1, T_2, T_3, T_4, T_6, T_8, T_9$
Milk	$T_3, T_5, T_6, T_7, T_8, T_9$
Cookie	$T_2, T_4$
Jam	$T_1, T_8$

item	Tidset
(B, B)	$T_1, T_4, T_8, T_9$
(B, m)	$T_5, T_7, T_8, T_9$
(B, C) $\times$	$T_4$
(B, J)	$T_1, T_8$
(Bu, m)	$T_3, T_6, T_8, T_9$
(Bu, C)	$T_2, T_4$
(Bu, J)	$T_1, T_8$
(m, C) $\times$	0
(m, J) $\times$	$T_8$
(C, J) $\times$	0

$K=3$ , MS = 2.

item	Tidset
B, Bu, M	$T_8, T_9$
B, Bu, J	$T_8, T_1$
B, Bu, C $\times$	$T_4$
B, m, J $\times$	$T_8$

$K=4$ .

item	Tidset
B, Bu, M, J	$T_8$

### Advantages:

- **Memory Requirements:** → ECLAT Algo. uses a Depth-first Search Approach, it uses less memory than Apriori.
- **Speed:** → The ECLAT Algo. is typically faster than Apriori Algo.
- **Number of Computations:** → The ECLAT Algo. does not involve the repeated scanning of the data to compute the individual support values.

## Decision Tree :-

Decision tree algorithm falls under the category of supervised learning. It can be used to solve both regression and classification problems. It is a tree that helps us in decision making purpose. The decision tree creates classification or regression models as a tree structure. Decision tree uses the tree representation to solve the problem. Decision tree contains ~~tree~~ 3 types of nodes.

① Root node:- It is the top most node in the tree. Data which is inside the node is known as attribute.

② Internal node:- Each internal node denotes a test on attribute. Nodes which are in between root node and leaf nodes are called as internal nodes.

③ Leaf node:- We call last nodes as leaf nodes. They represents output is class label.

- Advantages
- ① It does not require any domain knowledge.
  - ② Classification steps of decision tree are simple and fast.
  - ③ Missing values in data does not effect output.
  - ④ A decision tree model is automatic and does not require a standardization of data.

key factors:- Building a decision tree is all about discovery. attributes that return the highest data gain

### Entropy:-

Entropy refers to a common way to measure impurity. In the decision tree, it measures the impurity in dataset.

### Information Gain:-

Information Gain refers to the ~~define~~ decline in entropy after the dataset is split. It is also called as ~~an~~ entropy reduction.

### Example:-

Day.	Weather.	Temperature	Humidity	Wind.	play.
1	Sunny.	Hot	High	weak	No
2	Cloudy	Hot	High	Weak	Yes.
3	Cloudy	mild.	High	Strong	Yes
4	Rainy.	mild	High	Strong	No
5	Sunny.	mild.	Normal	Strong	Yes.
6	Rainy.	Cool	Normal	Strong	No
7	Rainy.	mild.	High	Strong	Yes
8	Sunny.	Hot	High	Weak	No
9	Cloudy	Hot	Normal	Strong	Yes
10	Rainy.	mild.	High	Strong	No

