## SNS COLLEGE OF TECHNOLOGY
### (AN AUTONOMOUS INSTITUTION)
### COIMBATORE – 35
### DEPARTMENT OF COMPUTER SIENCE AND ENGINEERING

UNIT IV          UNSUPERVISED LEARNING

Introduction - Association Rules – Apriori Algorithm - Clustering- K-means – EM Algorithm- Mixtures of Gaussians - Self-organizing Map - Principal Components, Curves and Surfaces – Independent Component Analysis. Case Study: Weather prediction.

## Association Rule Learning

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

The association rule learning is one of the very important concepts of machine learning and it is employed in **Market Basket analysis, Web usage mining, continuous production, etc.** Here 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 within a shelf or mostly nearby. Consider the below diagram:



Association rule learning can be divided into three types of algorithms:

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

We will understand these algorithms in later chapters.

**How does Association Rule Learning work?**

Association rule learning works on the concept of If and Else Statement, such as if A then B.



Here the If element is called **antecedent**, and then statement is called as **Consequent**. These types of relationships where we can find out some association or relation between two items is known *as single cardinality*. It is all about creating rules, and if the number of items increases, then cardinality also increases accordingly. So, to measure the associations between thousands of data items, there are several metrics. These metrics are given below:

- **Support**
- **Confidence**
- **Lift**

**Support**

Support is the frequency of A or how frequently an item appears in the dataset. It is defined as the fraction of the transaction T that contains the itemset X. If there are X datasets, then for transactions T, it can be written as:

$$Supp(X) = \frac{Freq(X)}{T}$$

**Confidence**

Confidence indicates how often the rule has been found to be true. Or how often the items X and Y occur together in the dataset when the occurrence of X is already given. It is the ratio of the transaction that contains X and Y to the number of records that contain X.

$$Confidence = \frac{Freq(X,Y)}{Freq(X)}$$

**Lift**

It is the strength of any rule, which can be defined as below formula:

$$Lift = \frac{Supp(X,Y)}{Supp(X) \times Supp(Y)}$$

It is the ratio of the observed support measure and expected support if X and Y are independent of each other. It has three possible values:

- If **Lift= 1**: The probability of occurrence of antecedent and consequent is **independent of each other.**
- **Lift>1**: It determines the degree to which the two itemsets are **dependent to each other**.
- **Lift<1**: It tells us that one item is a **substitute for other items**, which means one item has a negative effect on another.

**Types of Association Rule Lerning**

Association rule learning can be divided into three algorithms:

**Apriori Algorithm**

This algorithm uses frequent datasets to generate association rules. It is designed to work on the databases that contain transactions. This algorithm uses a **breadth-first search and Hash Tree** to calculate the **itemset efficiently**.

It is mainly used for **market basket analysis** and helps to understand the products that can be bought together. It can also be used in the healthcare field to find drug reactions for patients.

**Eclat Algorithm**

Eclat algorithm stands for **Equivalence Class Transformation**. This algorithm uses a **depth-first search technique to find frequent itemsets in a transaction database**. It performs **faster** execution than Apriori Algorithm.

**F-P Growth Algorithm**

The F-P growth algorithm stands for **Frequent Pattern**, and it is the improved version of the Apriori Algorithm. It represents the database in the form of a tree structure that is known as a frequent pattern or tree. The purpose of this frequent tree is to extract the most frequent patterns.

**Applications of Association Rule Learning**

It has various applications in machine learning and data mining. Below are some popular applications of association rule learning:

- **Market Basket Analysis:** It is one of the popular examples and applications of association rule mining. This technique is commonly used by big retailers to determine the association between items.
- **Medical Diagnosis:** With the help of association rules, patients can be cured easily, as it helps in identifying the probability of illness for a particular disease.
- **Protein Sequence:** The association rules help in determining the synthesis of artificial Proteins.
- It is also used for the **Catalog Design** and **Loss-leader Analysis** and many more other applications.

## Apriori Algorithm

Apriori algorithm refers to the algorithm which is used to calculate the association rules between objects. It means how two or more objects are related to one another. In other words, we can say that the apriori algorithm is an association rule leaning that analyzes that people who bought product A also bought product B.

The primary objective of the apriori algorithm is to create the association rule between different objects. The association rule describes how two or more objects are related to one another. Apriori algorithm is also called frequent pattern mining. Generally, you operate the Apriori algorithm on a database that consists of a huge number of transactions. Let's understand the apriori algorithm with the help of an example; suppose you go to Big Bazar and buy different products. It helps the customers buy their products with ease and increases the sales performance of the Big Bazar. In this tutorial, we will discuss the apriori algorithm with examples.

# Introduction

We take an example to understand the concept better. You must have noticed that the Pizza shop seller makes a pizza, soft drink, and breadstick combo together. He also offers a discount to their customers who buy these combos. Do you ever think why does he do so? He thinks that customers who buy pizza also buy soft drinks and breadsticks. However, by making combos, he makes it easy for the customers. At the same time, he also increases his sales performance.

Similarly, you go to Big Bazar, and you will find biscuits, chips, and Chocolate bundled together. It shows that the shopkeeper makes it comfortable for the customers to buy these products in the same place.

Current TimeÂ 4:42

/

DurationÂ 4:57

Â

X

The above two examples are the best examples of Association Rules in Data Mining

. It helps us to learn the concept of apriori algorithms.

## What is Apriori Algorithm?

Apriori algorithm refers to an algorithm that is used in mining frequent products sets and relevant association rules. Generally, the apriori algorithm operates on a database containing a huge number of transactions. For example, the items customers but at a Big Bazar.

Apriori algorithm helps the customers to buy their products with ease and increases the sales performance of the particular store.

## Components of Apriori algorithm

The given three components comprise the apriori algorithm.

1. Support
2. Confidence
3. Lift

Let's take an example to understand this concept.

We have already discussed above; you need a huge database containing a large no of transactions. Suppose you have 4000 customers transactions in a Big Bazar. You have to calculate the Support, Confidence, and Lift for two products, and you may say Biscuits and Chocolate. This is because customers frequently buy these two items together.

Out of 4000 transactions, 400 contain Biscuits, whereas 600 contain Chocolate, and these 600 transactions include a 200 that includes Biscuits and chocolates. Using this data, we will find out the support, confidence, and lift.

### Support

Support refers to the default popularity of any product. You find the support as a quotient of the division of the number of transactions comprising that product by the total number of transactions. Hence, we get

Support (Biscuits) = (Transactions relating biscuits) / (Total transactions)

= 400/4000 = 10 percent.

### Confidence

Confidence refers to the possibility that the customers bought both biscuits and chocolates together. So, you need to divide the number of transactions that comprise both biscuits and chocolates by the total number of transactions to get the confidence.

Hence,

Confidence = (Transactions relating both biscuits and Chocolate) / (Total transactions involving Biscuits)

= 200/400

= 50 percent.

It means that 50 percent of customers who bought biscuits bought chocolates also.

### Lift

Consider the above example; lift refers to the increase in the ratio of the sale of chocolates when you sell biscuits. The mathematical equations of lift are given below.

Lift = (Confidence (Biscuits - chocolates)/ (Support (Biscuits)

= 50/10 = 5

It means that the probability of people buying both biscuits and chocolates together is five times more than that of purchasing the biscuits alone. If the lift value is below one, it requires that the people are unlikely to buy both the items together. Larger the value, the better is the combination.

# How does the Apriori Algorithm work in Data Mining?

We will understand this algorithm with the help of an example

Consider a Big Bazar scenario where the product set is P = {Rice, Pulse, Oil, Milk, Apple}. The database comprises six transactions where 1 represents the presence of the product and 0 represents the absence of the product.

| Transaction ID | Rice | Pulse | Oil | Milk | Apple |
|---|---|---|---|---|---|
| t1 | 1 | 1 | 1 | 0 | 0 |
| t2 | 0 | 1 | 1 | 1 | 0 |
| t3 | 0 | 0 | 0 | 1 | 1 |
| t4 | 1 | 1 | 0 | 1 | 0 |
| t5 | 1 | 1 | 1 | 0 | 1 |
| t6 | 1 | 1 | 1 | 1 | 1 |

The Apriori Algorithm makes the given assumptions

- All subsets of a frequent itemset must be frequent.
- The subsets of an infrequent item set must be infrequent.
- Fix a threshold support level. In our case, we have fixed it at 50 percent.

## Step 1

Make a frequency table of all the products that appear in all the transactions. Now, short the frequency table to add only those products with a threshold support level of over 50 percent. We find the given frequency table.

**Product Frequency (Number of transactions)**

Rice (R)  4

Pulse(P) 5

Oil(O)    4

Milk(M) 4

The above table indicated the products frequently bought by the customers.

## Step 2

Create pairs of products such as RP, RO, RM, PO, PM, OM. You will get the given frequency table.

**Itemset Frequency (Number of transactions)**

RP    4

RO    3

RM    2

PO    4

PM    3

OM    2

**Step 3**

Implementing the same threshold support of 50 percent and consider the products that are more than 50 percent. In our case, it is more than 3

Thus, we get RP, RO, PO, and PM

**Step 4**

Now, look for a set of three products that the customers buy together. We get the given combination.

1. RP and RO give RPO
2. PO and PM give POM

**Step 5**

Calculate the frequency of the two itemsets, and you will get the given frequency table.

**Itemset Frequency (Number of transactions)**

RPO    4

POM    3

If you implement the threshold assumption, you can figure out that the customers' set of three products is RPO.

We have considered an easy example to discuss the apriori algorithm in data mining. In reality, you find thousands of such combinations.

# How to improve the efficiency of the Apriori Algorithm?

There are various methods used for the efficiency of the Apriori algorithm

**Hash-based itemset counting**

In hash-based itemset counting, you need to exclude the k-itemset whose equivalent hashing bucket count is least than the threshold is an infrequent itemset.

**Transaction Reduction**

In transaction reduction, a transaction not involving any frequent X itemset becomes not valuable in subsequent scans.

# Apriori Algorithm in data mining

We have already discussed an example of the apriori algorithm related to the frequent itemset generation. Apriori algorithm has many applications in data mining.

The primary requirements to find the association rules in data mining are given below.

**Use Brute Force**

Analyze all the rules and find the support and confidence levels for the individual rule. Afterward, eliminate the values which are less than the threshold support and confidence levels.

**The two-step approaches**

The two-step approach is a better option to find the associations rules than the Brute Force method.

**Step 1**

In this article, we have already discussed how to create the frequency table and calculate itemsets having a greater support value than that of the threshold support.

**Step 2**

To create association rules, you need to use a binary partition of the frequent itemsets. You need to choose the ones having the highest confidence levels.

In the above example, you can see that the RPO combination was the frequent itemset. Now, we find out all the rules using RPO.

RP-O, RO-P, PO-R, O-RP, P-RO, R-PO

You can see that there are six different combinations. Therefore, if you have n elements, there will be $2^n - 2$ candidate association rules.

# Advantages of Apriori Algorithm

- It is used to calculate large itemsets.
- Simple to understand and apply.

# Disadvantages of Apriori Algorithms

- Apriori algorithm is an expensive method to find support since the calculation has to pass through the whole database.
- Sometimes, you need a huge number of candidate rules, so it becomes computationally more expensive.