

# A Survey on Association Rule Mining using Apriori Algorithm

K.Saravana Kumar  
Ph.D Research Scholar  
NGM College,  
Pollachi.

R.Manicka Chezian  
Associate Professor  
NGM College,  
Pollachi.

## ABSTRACT

Association rule mining is the most important technique in the field of data mining. It aims at extracting interesting correlation, frequent pattern, association or casual structure among set of item in the transaction database or other data repositories. Association rule mining is used in various areas for example Banking, department stores etc. The paper surveys the most recent existing association rule mining techniques using apriori algorithm.

## Keywords

Association rule mining, Ant colony optimization, E-mail detection, Database Reverse Engineering.

## 1. INTRODUCTION

Association rule discovery from large databases is one of the tedious tasks in data mining. Suspicious e-mails are detected in two ways. Informative e-mails and alert e-mails. Informative e-mails give information about past criminal activities and Alert e-mails gives alert about future criminal activities using apriori algorithm. ACO algorithm gives the specific problem of minimizing the number of association rules [3]. Apriori algorithm uses transaction data set and uses a user interested support and confidence value and produce the association rule set. These association rule sets are discrete and continue therefore weak rule set are required to prune. Finally, conclude the paper.

## 2. RELATED WORKS

Determination of association rule mining, e-mail (emails about criminal activity) is suspected. Negative emotion words betray theory, a new person pronoun, in addition to simple words, the high-frequency words and special words were written in the body is characterized by deceptive e-mail writing preprocessed. Terms of apriori algorithm [1] is used to make. Data generated in the mail soon. It is used for automated analysis and evaluation to identify criminal activities and the announcement. Apriori algorithm for association rule mining, and all e-mail messages using the action verbs, past tense, using futures and evaluated. It's an action verb, such kind of emails in the future tense suffix and another with a message by e-mail if you are suspicious. Warning email to "“kill and bomb” future tense of words such as "“will and shall," which refers to such terms. Step number.

In order to classify the e-mail box, all HTML from the text element, header, body, etc removed, before the words are stop words tokenizing. After separation of the body, begins to move e-mail classification. Training data "Bomb/Blast/Kill " key and " will/may " they, important that the class information, e-mails, a move "attacked/terrorist" and tense "was" in them using apriori algorithm.. The training set apriori algorithm to find the e-mail database of words frequently used in the mining frequent item sets. Apriori algorithm for association rules and the rules used to set this item as follows.

Tense=past, Attack= Y, Bomb =Y ->Email = Suspicious informative Email.

Tense=future, Attack= Y, Bomb =Y-> Email =Suspicious alert Email.

Tense=future, Attack= N, Bomb =N->Email =Normal Email.  
An improved frequent pattern tree based on the technique named dynamic frequent pattern tree is proposed by Gyorodi [2]. The new method is efficiently applied on real world size database. A comparison between classical frequent pattern mining algorithms that are candidate set generation, test and without candidate generation is proposed in paper. Apriori algorithm, frequent pattern growth, dynamic frequent pattern growth are compared and presented together. Apriori algorithm in used to rule mining in huge transaction database and Apriori algorithm is a bottom up approach. Frequent pattern growth is used to novel, compact data structure, referred to as frequent pattern tree, fp tree based ones are partition based, divide and conquer methods.

Optimization of association rule mining and apriori algorithm Using Ant colony optimization [3].This paper is on Apriori algorithm and association rule mining to improved algorithm based on the Ant colony optimization algorithm. ACO was introduced by dorigo and has evolved significantly in the last few years. Many organizations have collected massive amount data. This data set is usually stored on storage database systems. Two major problems arise in the analysis of the information system. One is reducing unnecessary objects and attributes so as to get the minimum subset of attributes ensuring a good approximation of classes and an acceptable quality of classification. Another one is representing the information system as a decision table which shows

dependencies between the minimum subset of attributes and particular class numbers without redundancy. In Apriori algorithm, is working process explained in steps. Two step processes is used to find the frequent item set to join and prune. ACO algorithm was inspired from natural behavior of ant colonies. ACO is used to solve to numerous hard optimizations including the traveling salesman problem. ACO system contains two rules .One is local pheromone update rule, which is applied in constructing solution. Another one is global pheromone update rule which is applied in ant construction.ACO algorithm includes two more mechanisms, namely trail evaporation and optionally deamonactions.ACO algorithm is used for the specific problem of minimizing the number of association rules. Apriori algorithm uses transaction data set and uses a user interested support and confidence value then produces the association rule set. These association rule set is discrete and continues. Hence weak rule set are required to prune.

Association rule mining template is guided from XML document. XML is used in all areas of Internet application programming and is giving large amount of data encoded in XML [4]. With the continuous growth in XML data sources, the ability to extract knowledge from them for decision support becomes increasingly important and desirable. Due to the inherent flexibilities of XML, in both structure and semantics, mining knowledge in the XML Era is faced with more challenges than in the traditional structured world. This paper is a practical model for mining association rules from XML document.XML enabled association rule frame work that was introduced by Feng.XML AR frame works better than simple structured tree. The framework is flexible and powerful enough to represent simple and difficult structured association rules in XML document. But the best level of XML document model which is not yet implemented in has been proposed. The problem of mining XML association rules from the content of XML documents is based on user provided rule template. An implementation model is already introduced by Feng. Our practical model consists of the following steps Filtering, Generating Virtual Transactions, Finding Association Rules, Converting extracted rules of XML AR rules and Visualizing. Filtering and Generating virtual transactions are most important steps in this model Filtering step uses the XML-AR template and extracts only those parts of XML that are interesting for the user. The next step, defining a transaction context, based on tag nesting in XML document uses to generate virtual transactions that can be used as input format by association rule mining algorithms (e.g. Apriori). As an example, consider the problem of mining frequent associations among people who appear as coauthors, with our XML-AR template formulate. The statement is two parts(body and head) and the each part has 3-level xml fragment display generated virtual transactional based on XML AR template and xml fragment of dblp collection. The experimental results have found the coauthors and the keyword relationship, mining from xml document.

Database reverse engineering association rule mining [5] is based on. Classification of the document database system

could not be found in the poorly written, or even a particularly difficult task. The concept design of database reverse engineering to recover the database in an attempt to exit. Mining technique to detect the use of the concept plan proposing a strategy paper. They used the normal form. Classical database is a valuable asset to the organization. New technologies were developed in 1970 as a COBOL and the database, and mini - computer platforms, file systems using older programming languages. Even some of the databases are outdated concepts such as the hierarchical data model, designed, and maintained and adjusted to serve the current needs of modern companies was difficult for them. Classical databases, messaging systems and their structures are related to the contents of the move and changing. This document is no longer in my approach to system design, however, is hard to achieve; most companies in general are rising. Problems of migrating legacy databases to retrieve and database structures, database reverse engineering has been proposed. Reverse engineering process design and manufacturing process from the first to explore the objective devices and other hardware. This method is often used in World War II. Databases in real life mining association rules in applying the huge amount but it always creates a major problem. We select only the strong association rules, including the law of the filter element design. To find lost documents in database design and normalization policies and association mining techniques. Using formal analysis techniques for database reverse engineering process to improve design and builds. NoWARs (Normalization rules) means that a new analytical technique used to create a formal opinion and communication. Formal process and technology together to enable NoWars a means of union. Process a data mining association rules in a database and the dataset can be fed. Call to implement apriori algorithm to discover association rules and by NoWARs association rules in a database can save a form as a result. NoWARs normalization process used in the selection rules. Finally, the rules related to use program to create a table in the form of 3NF. Conceptual Schema and the power efficiency in the execution and analysis to filter out and recover. NoWARs maintains a database.Frequent Pattern - Classification: Generalized association rule mining based on tree structure. Into account [6] the use of specific information useful knowledge than ordinary flat association rules mining to detect this possibility. It employs a tree structure to compress the database Fp-line method, based on the paradigm proposing to mine development. There are two methods to the study of tree structure below and above are below. There are only a few studies have given way to common rules. Even some of the SQL queries and reports to the general association rule proposed parallel machine performance ratings. Apriori algorithm is a level-wise approach. A frequent pattern mining paradigm in the field of research has become a new trend. Often times, or to divide the so-called examples of successful development. Fp-growth algorithm is successful predecessor. Fp-Fp-growth trees that are often forms of devices to a data structure that collects all the required information. FP-tree traversal algorithm is the bottom line is above the upper fp-line and used to travel down and bottom up traversal methods

differ. Step three optimizations, such as the ancestors of transactions included in the overall filter, pre-computing the ancestors and the ancestors of an item is an item set, and pruning. With the apriori algorithm based on association rule mining algorithms based on common trees. The traversal of the tree structure below and above the bottom up traversal times worked construction. Efficient association rules mining properties by using the apriori algorithm. Apriori algorithm to obtain the means of association rules from the dataset. It's [7] the case of a large dataset is a time consuming procedure apriori algorithm is an efficient algorithm. Time changes many long-term activities to increase the number of paths apriori and the proposed database. Disadvantages and apriori algorithm apriori algorithm can improve performance; this paper describes the application properties. Customers who buy products at the beginning of an association rule mining is market basket analysis to find out how. Minimum support of all frequent items finding all association rules at  $\text{support} = 1$ . Find they considered the two-stage process. Enumeration of all frequent item sets the size of the search space is  $2^n$ . 2. To create strong rules. Used to create an association rule that satisfies any of the gate. Apriori and apriori algorithm using data mining tool and then running to write pseudo code. Association rule mining is the limit. ARM algorithm encounters a problem that does not return by the end user in a reasonable time. Activity in the presence and absence of an item is the only database that tells us a lot of shortcomings, and it is not efficient in the case of large dataset. ARM. The weight and size can be removed using such properties. Some of the limitations associated with large database apriori algorithm

Searches. Its easy to implement using apriori exists. Association rule mining as well as potential customers for commercial gain valuable information, much improved by the use of such properties. Association rule mining has wide applicability in many areas, efficient algorithm as it is a time consuming algorithm in case of large dataset [7]. With the time a number of changes are proposed in Apriori to enhance the performance in term of time and number of database passes. This paper illustrates the apriori algorithm disadvantages and utilization of attributes which can improve the efficiency of apriori algorithm. Association rule mining is initially used for Market Basket Analysis to find how items purchased by customers are related. The problem of finding association rules can be stated as follows: Given a database of sales transactions, it is desirable to discover the important associations among different items such the presence of some items in a transaction will imply the presence of other items in the same transaction . Discovering all association rules is considered as two-phase processes which are 1. Find all frequent item sets having minimum support. The search space to enumeration all frequent item sets is on the magnitude of  $2^n$ . 2. Generate strong rules. Any association that satisfies the threshold will be used to generate an association rule. The first phase in discovering all association rules is considered to be the most important one because it is time consuming due to the huge search space (the power set of the set of all items) and the second phase can be accomplished in a straight

forward manner. Then write the pseudo code for apriori and working of apriori algorithm using data mining tool. Limitation of association rule mining. The end user of ARM encounter problem as the algorithm does not return in a reasonable time. It is only tells the presence and absence of an item transactional database and it is not efficient in case of large dataset. ARM has a lot of disadvantages. These can be removed by using attributes like weight and quantity. Weight attribute will give user an estimate of how much quantity of item has been purchased by the customer, profit attribute will calculate the profit ratio and tell total amount of profit an item is giving to the customer. Apriori algorithm is associated with certain limitations of large database scans. Advantage of apriori is its easy implementation. Association rule mining efficiency can be improved by using attributes like profit, quantity which will give the valuable information to the customer as well as the business. Association rule mining has a wide range of applicability in many areas

### **3. DATABASE REVERSE ENGINEERING**

Legacy databases are important assets to organizations. However these databases were mostly developed using 1970's technologies and hence is difficult to adapt to the current needs of modern companies. Maintenance of these legacy databases is very difficult mainly, due to the lack of proper system documentation. So in-order to ease the process of maintaining these legacy databases the concept of Database reverse Engineering arose. Database reverse engineering solves the problem of recovering database structure and migrating legacy database. Database reverse engineering refers to the process of discovering the source code as well as the system design of the given software. To reverse engineer a database is, to discover the underlying relations between various records and tables such as the keys, functional dependencies and integrity constraints. Typically this means applying association mining to a database, but in real life databases this leads to the generation of a tremendous amount of association rules. Apply the normalization principles and association mining techniques to discover the missing database design. To control the number of associations another more efficient technique named 'NoWARs' (Normalization With association Rules), which uses the Apriori algorithm, was developed. The NoWARs technique user the concept of normalization and the concept of association, it consists of two important steps: finding the rules of association and then normalizing those rules. The algorithm NoWARs starts when used enter query to define dataset to normalize. Then NoWARs will find the association rules by calling Apriori algorithm and save resulting a form of association rules in the database. Then NoWARs will select some rules to use in normalization process. Finally, use the selected rules to generate the 3NF in relational schema form. NoWARs are normalizing the database tables. In the normalization process, NoWARs uses only 100% confidence association rules with any support values. Recent work in database reverse engineering has not concentrated on a broad

objective of system migration; instead focus is made on a particular issue of semantic understanding. Reverse engineering method to discover inter-relational constraints and inheritances embedded in a relational database.

#### 4. CONCLUSION

Association rule mining has a wide range of applicability such as market basket analysis, suspicious e-mail detection, library management and many areas. The conventional algorithm of association rules discovery proceeds in two steps. All frequent item sets are found in the first step. The frequent item set is the item set that is included in at least minimum support transactions. The association rules with the confidence at least minimum confident are generated in the second step. In this paper, we surveyed the list of existing association rule mining techniques using apriori algorithm.

#### 5. REFERENCES

- [1] S.Appavu alias Balamurugan, Aravind , Athiappan, Barathiraja, Muthu Pandian and Dr.R.Rajaram, "Association rule mining for suspicious Email detection:A data mining approach" 11-4244-11330-3/07/\$25.00 02007 IEEE.
- [2] Cornelia Györödi, Robert Györödi, T. Cofeey & S. Holban – "Mining association rules using Dynamic FP-trees" – in proceedings of The Irish Signal and Systems Conference, University of Limerick, Limerick, Ireland, 30th June-2nd July 2003, ISBN 0-9542973-1-8, pag. 76-82.
- [3] Badri patel ,Vijay K Chaudahri,Rajneesh K Karan,YK Rana --"Optimization of association rule mining apriori algorithm using Ant Conoly optimization" International Journal of Soft Computing and Engineering (IJSCE) ISSN: 2231-2307, Volume-1, Issue-1, March 2011
- [4] Rahman AliMohammadzadeh, Sadegh Aoltan, Masoud Rahgozar –"Template Guided Association Rule Mining from XML Documents" WWW 2006, May 23–26, 2006, Edinburgh, Scotland.ACM 1-59593-323-9/06/0005.
- [5] Nattapon Pannurat, Nittaya Kerdprasop, Kittisak Kerdprasop –"Database Reverse Engineering based on Association Rule Mining" IJCSI International Journal of Computer Science Issues, Vol. 7, Issue 2, No 3, March 2010 ISSN (Online): 1694-0784ISSN (Print): 1694-0814
- [6] Iko Pramudiono and masaru Kitsuregawa - "FP-tax :Tree structure Generalized Association Rule Mining" DMKD '04, June 13, 2004, Paris, France Copyright 2004 ACM ISBN 158113908X/04/06 ...\$5.00.
- [7] Mamta Dhanda, sonali Guglani and gaurav gupta – "Mining Efficient Association Rules through Apriori Algorithm Using Attributes" IJCST Vol. 2, Issue 3, September 2011 I S S N : 2 2 2 9 - 4 3 3 3 ( P r i n t ) | I S S N : 0 9 7 6 - 8 4 9 1 (O n l i n e )
- [8] J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann Publishers, San Francisco, USA, 2001, ISBN 1558604898.
- [9] L. Cristofor, "Mining Rules in Single-Table and Multiple-Table Databases",PhD thesis, University of Massachusetts, 2002.
- [10] Ashoka Savasere, Edward Omiencinski, and Shamkant B. Navathe. "AnEfficient Algorithm for Mining Association Rules in Large Databases." InProceedings of the 21st International Conference on Very Large Databases, pag432 - 444, 1995.