# Method for extracting Frequent Pattern Using Transposition of Database

Birla Sunderlal

**Abstract:** The purpose of the Apriori Algorithm is to find frequent itemsets between different transaction sets of data.Apriori is a classical algorithm for frequent patterns extraction. An approach implemented in Transposed database then result is very fast. Apriori is work  to operate on databases containing transactions. The main  aim of this research is to improve the performance of the conventional Apriori algorithm that extracts frequent patterns for binary transaction dataset. Recently, different works proposed a new way to mine frequent patterns in transposed databases where a database with thousands of attributes but only tens of objects. In this case, mining the transposed database runs through a smaller search space. This work systematically explores the search space of frequent patterns mining and represent database in transposed form. This paper proposed an algorithm for mining frequent patterns which are based on Apriori algorithm and used space reduced longest common sequence (LCS) which makes apriori algorithm space efficient. Space complexity for Proposed algorithm is $O(n)$ while the Dynamic Approach like Longest Common Subsequence space complexity is $O(n^2)$ memory for given items in dataset.

**Keywords:** Transposition of database, Apriori algorithm, LCS, Frequent itemset, Data mining, Space complexity.
.

## I. INTRODUCTION

Apriori is a classic algorithm for learning association rules. Apriori is designed to operate on databases containing transactions. Other algorithms are designed for finding association rules in data having no transactions, or having no timestamps. Other algorithms are designed for finding association rules in data having no transactions, or having no timestamps.  The purpose of the Apriori Algorithm is to find associations between different sets of data. It is sometimes referred to as "Market Basket Analysis". Each set of data has a number of items and is called a transaction.  The output of Apriori is sets of rules that tell us how often items are contained in sets of data. Frequent item set mining and association rule induction are powerful methods for so-called market basket analysis, which aims at finding regularities in the shopping behavior of customers of supermarkets, mail-order companies, online shops etc. With the induction of frequent item sets and association rules one tries to find sets of products that are frequently bought together, so that from the presence of certain products in a shopping cart one can infer (with a high Probability) that certain other products are present. Such information, especially if expressed in the form of rules, can often be used to increase the number of items sold, for instance, by appropriately arranging the products on the shelves of a supermarket or on the pages of a mail-order catalog (they may, for example, be placed adjacent to each other in order to invite even more customers to buy them together) or by directly suggesting items to a customer, which may be of interest for him/her.

## II. 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 pre-processed. 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 associetion 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 co-authors, 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 atuvallatu

1. Find they considered the two-stage process.  Enumeration of all frequent item sets the size of the search space is 2n.
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

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.

### III. TRANSPOSITION OF DATABASE

To avoid confusion between rows (or columns) of the original database and rows (columns) of the "transposed" database, we define a database as a relation between original and transposed representations of a database in Table-1. The attributes are A = {a1, a2, a3, a4} and the objects are O = {o1, o2, o3}. We use a string notation for object sets or itemsets, e.g., a1a3a4 denotes the itemset {a1, a3, a4} and o2o3 denotes the object set {o2, o3}. This dataset is used in all the examples between two sets: a set of attributes and a set of objects.

TABLE I.   TRANSPOSITION OF DATABASE DATABASE D

| TRANSPOSED | | DATABASE D$^T$ | |
|---|---|---|---|
| Object | Attribute Pattern | Attribute | Object Pattern |
| O1 | a1a2a3 | a1 | O1O2 |
| O2 | a1a2a3 | a2 | O1O2O3 |
| O3 | a2a3a4 | a3 | O1O2O3 |
| | | a4 | O3 |

### A.  DYNAMIC APPROACH

The longest common subsequence problem is one of the common problems which can be solved efficiently using dynamic programming. "The Longest common subsequence problem is, we are given two sequences X=<x1,x2----------xn> and Y=<y1,y2---------ym> and wish to find a maximum length common subsequence of X and Y" for example : if X=<A,B,C,B,D,A,B>  and  Y=<B,D,C,A,B,A>  then  The sequence <B, C, B, A> longest common subsequence. Let us define C [i, j] to be the length of an LCS of the sequences xi and yj. If either i=0 or j=0, one of the sequence has length 0, so the LCS has length 0. The Optimal substructure of the LCS Problem gives the recursive formula in fig.1

$$C(i, j) = \begin{cases} 0 & if\ i = 0\ or\ j = 0 \\ C(i-1, j-1) + 1 & if\ i, j > 0\ and\ xi = yj \\ \max(c(i, j-1), c(i-1, j)) & if\ i, j > 0\ and\ xi \neq yj \end{cases}$$

Fig. 1 Longest Common Subsequence Recursive FormulaProposed Reduced space complexity

Back bone of this algorithm is finding LCS from dataset which is perform by dynamic algorithm in above work we proposed Heuristic space efficient algorithm But it can also be thought of as a way of computing the entries in the array L. The recursive algorithm controls what order we fill them in, but we'd get the same results if we filled them in in some other order. We might as well use something simpler, like a nested loop, that visits the array systematically. The only thing we have to worry about is that when we fill in a cell L[i,j], we need to already know the values it depends on, namely in this case L[i+1,j], L[i,j+1], and L[i+1,j+1]. For this reason we'll traverse the array backwards, from the last row working up to the first and from the last column working up to the first. This is iterative (because it uses nested loops instead of recursion) or bottom up (because the order we fill in the array is from smaller simpler subproblems to bigger more complicated ones).

### A.  Iterative LCS:

```
int lcs_length(char * A, char * B)
{
    allocate storage for array L;
    for (i = m; i >= 0; i--)
        for (j = n; j >= 0; j--)
        {
            if (A[i] == '\0' || B[j] == '\0') L[i,j] = 0;
            else if (A[i] == B[j]) L[i,j] = 1 + L[i+1, j+1];
            else L[i,j] = max(L[i+1, j], L[i, j+1]);
        }
    return L[0,0];
}
```

Advantages of this method include the fact that iteration is usually faster than recursion, we don't need to initialize the matrix to all -1's, and we save three if statements per iteration

since we don't need to test whether L[i,j], L[i+1,j], and L[i,j+1] have already been computed (we know in advance that the answers will be no, yes, and yes). One disadvantage over memoizing is that this fills in the entire array even when it might be possible to solve the problem by looking at only a fraction of the array's cells

One disadvantage of the dynamic programming methods we've described, compared to the original recursion, is that they use a lot of space: O(mn) for the array L (the recursion only uses O(n+m)). But the iterative version can be easily modified to use less space -- the observation is that once we've computed row i of array L, we no longer need the values in row i+1.

*B. Space-efficient LCS:*

```
int lcs_length(char * A, char * B)
{
    allocate storage for one-dimensional arrays X and Y
    for (i = m; i >= 0; i--)
    {
        for (j = n; j >= 0; j--)
        {
            if (A[i] == '\0' || B[j] == '\0') X[j] = 0;
            else if (A[i] == B[j]) X[j] = 1 + Y[j+1];
            else X[j] = max(Y[j], X[j+1]);
        }
        Y = X;
    }
    return X[0];
}
```

This takes roughly the same amount of time as before, O(mn) or O(n2) it uses a little more time to copy X into Y but this only increases the time by a constant (and can be avoided with some more care). The space is either O(m) or O(n), whichever is smaller (switch the two strings if necessary so there are more rows than columns). Unfortunately, this solution does not leave you with enough information to find the subsequence itself, just its length.

## IV. EXPLANATION WITH EXAMPLE WHICH SUPPORT THE ARGUMENTS

Study the following transaction database. A={A1, A2, A3, A4, A5, A6, A7, A8, A9>, Assume  σ=20%, Since T contains 15 records, it means that an itemset that is supported by at least three transactions is a frequent set

TABLE II.  GIVEN DATASE T

| S. No. | A1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 | A9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 |
| 2 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 3 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 4 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| 6 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |
| 8 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 10 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| 11 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 |
| 13 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 |
| 14 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| 15 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |

TABLE III.  TRANSPOSITION OF DATABASE WITH TRANSACTION ID

| Item Id | Transaction id String | Count |
|---|---|---|
| 1 | 1,14 | 2 |
| 2 | 2,4,6,7,13,15 | 6 |
| 3 | 4,6,10,11,14,15 | 6 |
| 4 | 2,3,6,13 | 4 |
| 5 | 1,3,5,8,10,11,12,14 | 8 |
| 6 | 1,5,7,12,13 | 5 |
| 7 | 3,5,7,10,11,13,14 | 7 |
| 8 | 1,2,9,12 | 4 |
| 9 | 7,15 | 2 |

Now Apply Algorithm
1) Pass 1

TABLE IV.  FREQUENT ITEM SET

| Item Id | Transaction id String | Count |
|---|---|---|
| 2 | 2,4,6,7,13,15 | 6 |
| 3 | 4,6,10,11,14,15 | 6 |
| 4 | 2,3,6,13 | 4 |
| 5 | 1,3,5,8,10,11,12,14 | 8 |
| 6 | 1,5,7,12,13 | 5 |
| 7 | 3,5,7,10,11,13,14 | 7 |
| 8 | 1,2,9,12 | 4 |

L1 :={ 2, 3, 4, 5, 6, 7, 8}
2) Pass 2
Generate candidate for k=2
C2:={{2,3},{2,4},{2,5},{2,6},{2,7},{2,8},{3,4},{3,5},{3,6},{3,7},{3,8},{4,5},{4,6},{4,7},{4,8},{5,6},{5,7},{5,8}, {6, 7}, {6, 8}, {7, 8}}
After Apply DAPS Algorithm

TABLE V.  RESULT OF DAPS ALGORITHM

| Item Id | Transaction id String | Count |
|---|---|---|
| 2,3 | 4, 6,15 | 3 |
| 2,4 | 2, 6,13 | 3 |
| 3,5 | 10,11,14 | 3 |
| 3,7 | 10,11,14 | 3 |
| 5,6 | 1,5,12 | 3 |
| 5,7 | 3,5,10,11,14 | 5 |
| 6,7 | 5,7,13 | 3 |

3) Pass 3
Generate Candidate for k: = 3
C3:={{2,3,4},{3,5,7},{5,6,7}}

After Apply DAPS Algorithm

| Item Id | Transaction id String | Count |
|---|---|---|
| 3,5,7 | 10,11,14 | 3 |
| 5,6,7 | 5 | 1 |

L3 :={( 3, 5, 7)}
L: =L1UL2UL3

## V.  CONCLUSION

Determining frequent objects (item sets, episodes, sequential Patterns) is one of the most important fields of data mining. It is well known that the way candidates are defined has great effect on running time and memory need, and this is the reason for the large number of algorithms. We presented a new research trend on frequent pattern mining is expecting transpose representation to relieve current methods from the traditional bottleneck, providing scalability to massive Data sets and improving response time. In order to mine patterns in databases with more columns than rows, we proposed a complete framework for the transposition: we gave the item set

in the transposed database of the transposition of much classical transaction ID. Then we gave a strategy to use this framework to mine all the itemset satisfying. We used space reduced approach which is better than tradition approach for finding longest common subsequence with efficient memory space for large dataset set because of its linear space complexity. Our algorithm takes O(n) space in place of tradition dynamic algorithm takes O(n2) Its also efficient for large value of n whether traditional algorithm is not suitable for large value of n. due to less memory use for LCS the overall memory use expected to reduce a lot.

## REFERENCES

[1]   S.Appavu alias Balamurugan, Aravind , Athiapppan, Barathiraja, Muthu Pandian and Dr.R.Rajaram," Association rule mining for suspicious Email detection:A data mining approach" l1-4244-1l330-3/07/$25.OO 02007 IEEE.

[2]   Cornelia Győrödi, Robert Győrödi, T. Cofeey & S. Holban – "Mining association rules using Dynamic FP-trees" –în 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 (On 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.

[11]  B. Jeudy and F. Rioult, Database transposition for constrained closed pattern mining, in: Proceedings of Third International Workshop on Knowledge Discovery in Inductive Databases (KDID) co-located with ECML/PKDD, 2004.

[12]  R. Agrawal, R. Srikant, Fast algorithms for mining association rules,In Proceedings of the 20th International Conference on Very Large Data Bases, 1994, pp. 487–499.

[13]  J. Han, Research challenges for data mining in science and engineering. In NGDM 2007.

[14]  R. Agrawal, R. Srikant, Mining sequential patterns, In Proceedings of the 11th International Conference on Data Engineering, 1995, pp. 3

[15]  A fast APRIORI implementation Ferenc Bodon, Informatics Laboratory, Computer and Automation Research Institute, Hungarian Academy of Sciences H-1111 Budapest, L´agym´anyosi u. 11, Hungary

[16]  B. Goethals. Survey on frequent pattern mining. Technical report, Helsinki Institute for Information Technology,03.

[17]  R. Agrawal and R. Srikant. Fast algorithms for mining association rules. The International Conference on Very Large Databases, pages 487–499, 1994.