

An Improved Algorithm of Graph and Clustering based Association Rule Mining (GCBARM) in Discovering of Frequent Itemsets

N. Balaji Raja*
G. Balakrishnan**

Abstract

The key process of Association Rule mining algorithms is discovering frequent itemsets. The new proposed numerous algorithms are an important research topic in the field of data mining. This paper proposed two concepts. First, in huge amount of datasets, basic need of valuable analysis is to finding the relational and geometric characteristics of the underlying entities are represent their relationships with vertices and edges, this method provide to represent such datasets. Second is the proposed algorithm, which based on graph and clustering based mining association rules. This proposed algorithm is named Graph and Clustering Based Association Rule Mining (GCBARM). Scanning is the important task of frequent itemsets, the GCBARM algorithm scans the database of transaction only once to generate a cluster table and then clusters the transactions into cluster according to their length. The GCBARM algorithm is used to find frequent itemsets and will be extracted directly by scanning the cluster table. This method reduces memory requirement and time to retrieve the datasets and hence it is scalable for any large size of the database.

Keywords: Association rule mining, Data Mining, Relational, cluster, graph, GCBARM

1. Introduction

With the rapid development in size and number of available databases in commercial, organizational, administrative and other applications [1, 19], it is the urgent need for new technologies and automated tools to change this wealth of data resources into useful information. The most challenging in database mining is developing fast and efficient algorithms that can deal with large volume of data because several data mining algorithm computation is used to solve diverse data mining problem. They are mainly classified as associations, classifications, sequential patterns and clustering [2].

1.1. Association Rule

Association rule mining is a very important research topic in the data mining field, it's problem in large databases is to generate all association rules [17, 20], of the form $X \Rightarrow Y$, that will produce strong association rules which satisfy both minimum support degree (min_sup) and the minimum confidence degree (min_conf) greater then the user defined minimum support and minimum confidence [2, 3, 4].

Definition 1: let $X = \{x_1, x_2, \dots, x_n\}$ be a set of items, then $D = \{ \langle T_{id}, T \rangle | T \subseteq X \}$ is a transaction database, where T_{id} is an identifier which be associated with each transaction.

*Department of Computer Applications, Tamilnadu, India

**Department of Computer Science and Engineering, Tamilnadu, India

Definition 2: Let, $A \subseteq X, B \subseteq X, \text{ and } A \cap B = \phi$ we call this $A \Rightarrow B$ as association rule.

Most of the algorithm generally is executed in two steps. First, to finding all sets of items that have support above the given minimum, and then generating the desired rules from these item sets. The apriori algorithm is the same as the association rule. The first step is to find all frequent itemsets, the next is to generate strong association rules from frequent item sets during pruning [18].

1.2. Apriori Algorithm

The apriori algorithm is a fast algorithm for mining association rules and is based on [5,20] algorithms for mining association rules, the problem of this algorithm is number of data scans n , where n is the size of large nonempty itemset and number of discovering rules is huge while most of the rules are not interesting. Therefore several improved algorithms were proposed after apriori for efficiency and scalability.

This paper introduces an algorithm GCBARM, which is basically different from previous algorithms; the remaining paper is organized as follows: related work is presented in Section 2, Section 3 gives the details of GCBARM algorithm with example, Section 4 show the experimental result and finally the conclusion is given in Section 5.

2. Literature Review on Association Rule and Graph Database

Existing studies in data mining have presented efficient algorithm for discovering association rules. But the main drawback of the first algorithm is the need to do multiple passes over the datasets to generate frequent itemsets. The apriori association rule algorithm proposed by Agrawal and Srikant [6] can discover meaningful itemsets, but a large number of the candidate itemsets are generated from single itemsets and level by level in the process of creating association rules. Performance is severely affected because the database is scanned repeatedly to each candidate itemset with the database.

FP-Growth [7] out performs all candidate set generation and test algorithms as it mines frequent patterns without candidate generation. The main problem is no common prefixes within the data items.

Sample algorithm reduces the scanning to the datasets, it's scans single scan, but wastes considerable time on candidate itemsets [8].

The column wise apriori algorithm [9] and the tree based association rule algorithm [10] transformed the storage structure of the data, to reduce the time of scans of the transaction database.

The partition algorithm to improve efficiency and reduce the database scans, but still wasted scans on infrequent candidate itemset [11].

The primitive association rule mining is mining that describe the association among items of the transaction in the database. A

uniform frame work was framed to perform association rule of association rules [12].

In the generalized association patterns, one can add all ancestors for each items from concept hierarchy and then apply the algorithm on the extended transactions [13].

The multiple level association rules are discovered from a large database of customer transactions in which all items are described by a set of relevant attributes. Each attribute represents a certain concept and these relevant attributes form a set of multiple level concepts [5].

For frequent patterns in this new graph model which we call taxonomy superimposed graphs, there may be many patterns that are implied by the generalization and specialization hierarchy of the associated node label taxonomy [14]

Graph databases are able to represent as graphs of any kind of information, where naturally accommodated changes in data can be possible [15, 16].

The disadvantage of these algorithms is

- Number of reads the database transaction n time data scans where n is the size of large nonempty itemset,
- It is an incompetent as it requires wastage memory.
- The number of discovered rules is huge while most of them are non interesting.

The proposal method GCBARM is a consequence to overcome the above said drawbacks.

3. Proposed Methodology

This is a chance to establish that cluster based algorithms are still available by providing a better graph data structure which is used to simplify the process of generating frequent k itemsets, where $k \geq 2$. The Graph and Clustering Based Association Rule Mining (GCBARM) algorithm scans the database of transaction only once, which overcome the drawbacks of the previous algorithms.

The process of building the graph is given in sequential numbers, this simplified pre-process is taken in to consideration as an important action before applying our proposed algorithm. The GCBARM algorithm scans the database of transaction only once to generate a cluster table as a two dimensional array where the rows represent transactions' TIDs and the columns represent items. The presence and absence of an item in a transaction indicates that 1 and 0 is content of the table. Here providing example of Adult datasets from UCI data collections, the dataset contains Age, work class, fnlwtg, education, education-num, marital-status, occupation, relationship, race, sex, capital-gain, capital-loss, hours-per-week, native-country and Prediction task is to determine whether a person makes over 50K a year.

Initially an original database consists of three relations, but it's pre-processed and shown in table 1. Then the three relations are converted into equivalent graph database (GDB) shown in figure 1, using this processing to identify the attributes for rule mining. The database consists of three relations, named Adult Personnel Detail, Adult Employment Detail and Qualification. The data

structure of the relations is

- Adult Personnel Detail (Fnlwgt, Age, marital-status, relationship, race, sex, native-country)
- Adult Employment Detail (Work class, fnlwgt, occupation, capital-gain, capital-loss, hours per week, prediction task 50k)
- Qualification (Fnlwgt, Education, education num)

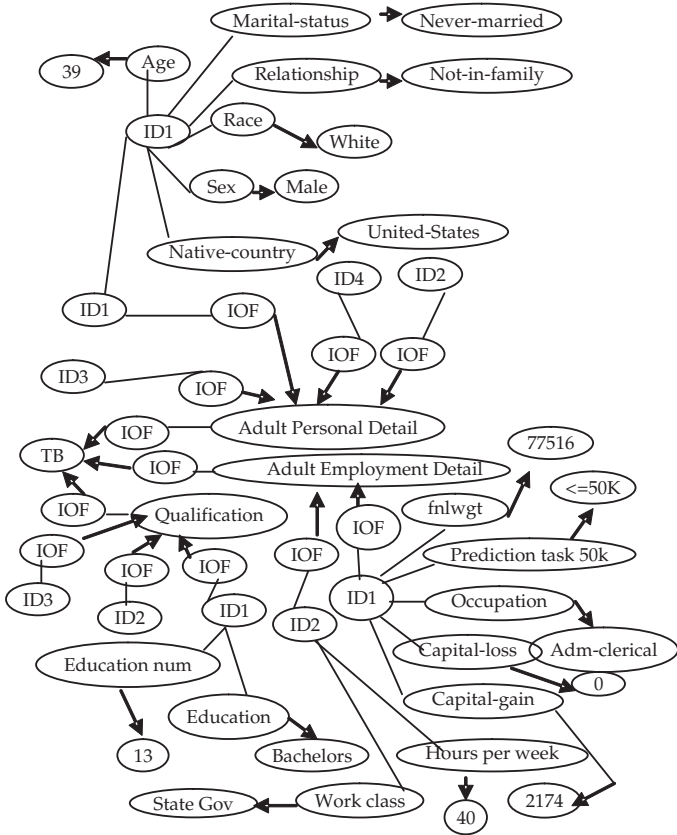


Fig. 1: Mapping of above relation into graph database

After pre-processing the set minimum support threshold is 50%. There are 20 transactions and five different items named Age, Work class, Marital status, Native country and prediction task in the database. Most of the rule mining algorithms are in lexicographical order. An example of transaction database is shown in Table1. The items name of the transaction database is used rather than numbers to deal with some worst cases. The requirement of numbering system is to first scan the database to identify the length of each transaction, that means length of the numbers to the items in a transaction, and at the same time assigning the numbers to the items: Number 1 is assigned as item Age, Number 2 is assigned as item work class, Number 3 is assigned as marital status, Number 4 is assigned as native country, Number 5 is assigned as prediction task. This conversion process help us in both constructing the cluster table and building the graph, this process help avoid the need to rescan the transaction database. Next move is the clustering table that can easily reside in the main memory.

In this example, the maximum transaction length is five, there will be at most five clusters, the total number of clusters is five as shown in table 2, The presence and absence of an item in a

Table 1: An example of transaction database

TID	Age	work class	marital-status	native-country	P_task
T ₁	39	State-gov	Never-married	United-States	<=50K
T ₂	50	Self-emp-not-inc	Married-civ-spouse	United-States	<=50K
T ₃	38	Private	Divorced	United-States	<=50K
T ₄	53	Private	Married-civ-spouse	United-States	<=50K
T ₅	28	Private	Married-civ-spouse	Cuba	<=50K
T ₆	37	Private	Married-civ-spouse	United-States	<=50K
T ₇	49	Private	Married-spouse-absent	Jamaica	<=50K
T ₈	52	Self-emp-not-inc	Married-civ-spouse	United-States	>50K
T ₉	31	Private	Never-married	United-States	>50K
T ₁₀	42	Private	Married-civ-spouse	United-States	>50K
T ₁₁	37	Private	Married-civ-spouse	United-States	>50K
T ₁₂	30	State-gov	Married-civ-spouse	India	>50K
T ₁₃	23	Private	Never-married	United-States	<=50K
T ₁₄	54	Private	Separated	United-States	<=50K
T ₁₅	35	Federal-gov	Married-civ-spouse	United-States	<=50K
T ₁₆	43	Private	Married-civ-spouse	United-States	<=50K
T ₁₇	59	Private	Divorced	United-States	<=50K
T ₁₈	56	Local-gov	Married-civ-spouse	United-States	<=50K
T ₁₉	19	Private	Never-married	United-States	>50K
T ₂₀	23	Local-gov	Never-married	United-States	<=50K

transaction is denoted by 1 and 0 in content of the table. After that, the bit vector for each item will be ready and it is an easy process to determine the frequent 1 itemsets by counting the number of 1s in each transaction, the minimum support threshold is not less than counting the number of 1s, but it is considered as a frequent itemset and then building the graph.

By counting the number of 1s in each bit vector is determined the support for each candidate itemset of length 1, as follows: support ({1}) = 55%, support ({2}) = 25%, support ({3}) = 55%, support ({4}) = 85%, support ({5}) = 65%. Thus the frequent 1 itemsets are :{{1}, {3}, {4}, {5}} as their supports are not less than 50%.

The second step starts by reordering frequent 1 itemsets by providing each one with a sequential number to facilitate the process of constructing the graph, which is making logical() and operation between each pair of consecutive frequent 1 itemsets<itemi,itemj> | i<j if the number of 1s in

Table 2: The cluster table form the database in table 1

Item No. TID	1	2	3	4	5
T7	0	0	0	0	1
T8	0	0	1	1	0
T9	1	0	0	1	0
T10	0	0	1	1	0
T14	0	0	0	1	1
T17	0	0	0	1	1
T19	1	0	0	1	0
T2	0	0	1	1	1
T3	1	0	0	1	1
T4	0	0	1	1	1
T5	1	0	1	0	1
T11	1	0	1	1	0
T12	1	1	1	0	0
T13	1	0	0	1	1
T16	0	0	1	1	1
T18	0	1	1	1	0
T1	1	1	0	1	1
T6	1	0	1	1	1
T20	1	1	0	1	1
T15	1	1	1	1	1

The bit vectors for the items are:

BV1=00100010101111001111
 BV2=00000000000010011011
 BV3=01010001011110110101
 BV4=01111111110101111111
 BV5=10001101111001101111

the result is greater than or equal to minimum support threshold, a edge is directed to drawn from item_i to item_j, this process is repeated for all frequent 1 itemsets. The simple directed graph to display frequent k - itemsets, k>=2 is shown in figure 2, and by assigning 25% as a new value to the minimum support threshold, the frequent 2 itemsets will be: {{1,3}, {1,4}, {3,5}} as shown in table 3 and the graph is constructed by drawing an edge between each pair of frequent items, as shown in figure 2. By counting the number of 1s in each bit vector is determined the support for each candidate itemset of length 2, as follows: support ({1,3}) = 25%, support ({1,4}) = 45%, support ({3,5}) = 30%. Thus the frequent 2 itemsets are :{{1,3}, {1,4}, {3,5}} as their supports are equal and above 25%.

Table 3: frequent itemsets 2 from frequent itemsets 1

Item No. TID	{1,3}	{1,4}	{3,5}
T7	0	0	0
T8	0	0	0
T10	0	0	0
T14	0	0	0
T17	0	0	0
T18	0	0	0
T9	0	1	0
T19	0	1	0

T2	0	0	1
T3	0	1	0
T4	0	0	1
T12	1	0	0
T13	0	1	0
T16	0	0	1
T1	0	1	0
T20	0	1	0
T5	1	0	1
T11	1	1	0
T6	1	1	1
T15	1	1	1

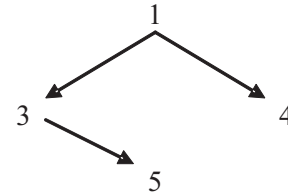


Fig. 2: A simple directed graph to display frequent k-itemsets, k>=2

The traverse of graph as if their path is to determine frequent 3 itemsets among three nodes{i,j} and {j,k} then the set {i,j,k} will be frequent 3 itemsets. Here, in this example, {{1,3,5}} is the only frequent 3 itemsets. As there are no extra edges, by assigning 12.5% as a new value to the minimum support threshold, the frequent 3 itemsets will be: {1,3,5}as shown in table 4 and the support for each candidate itemset of length 3, as follows: support ({1,3,5}) = 15%. Thus the frequent 3 itemsets are :{1,3,5} as their supports are equal and above 25%. Finally the algorithm terminates

Table 4: frequent itemsets 3 from frequent itemsets 1

Item No. TID	{1,3,5}
T7	0
T8	0
T9	0
T10	0
T14	0
T17	0
T19	0
T2	0
T3	0
T4	0
T11	0
T12	0
T13	0
T16	0
T18	0
T1	0
T20	0
T5	1
T6	1
T15	1

In this regard, as the database contains hundreds and thousands of transactions database and different items, constructing only one graph is not suitable for this work: So construct different graphs for each cluster and find from this graph all frequent itemsets, then combine the subsets of frequent itemsets together to get the whole set of frequent itemsets, and this technique is scalable with all transactions databases of different sizes.

4. Experimental Results

In order to appraise the performance of the proposed technique, which was implemented the GCBARM algorithm, along with Apriori algorithm using java programming language. The test databases are collected from UCI standard datasets available to evaluate rule mining algorithms.

Apriori and GCBARM are execute both algorithm at various values of minimum support thresholds, Figure 3 shows the average execution time (seconds) to generate all frequent itemsets using GCBARM and Apriori. The experimental results in figure 3 show the better performance of GCBARM algorithm than Apriori in terms of execution time.

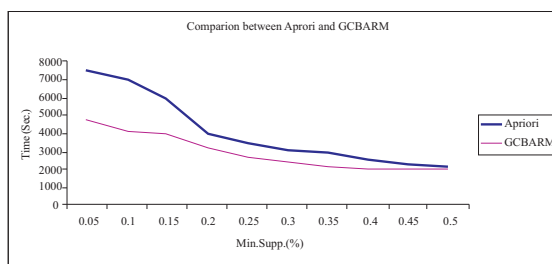


Fig. 3: Comparison between Apriori and GCBARM

5. Conclusion

This paper proposes the use of the graph database for pre-processing, after the hole transaction database is divided into partitions of variable sizes. Each cluster is considered one at a time by loading the first cluster into memory and calculating large itemsets and the corresponding support counts. Then the second cluster is considered similarly and the cumulative support count is calculated for the cumulative large itemsets. This process is continued for the entire set of clusters and finally the whole large itemsets and the corresponding cumulative support counts. This approach reduces main memory requirement since it considers only a small cluster at a time and hence it is scalable for any large size of the database.

6. References

1. Willi Klogsen and Jan M. Zytkow, "Hand Book of Data Mining and Knowledge Discovery", Oxford University Press, 2002.
2. Jiawei Han, Micheline Kamber, "Data Mining Concepts and Techniques", Morgan Kauffman Publishers, San Francisco 2001.
3. Rakesh agrawal, Tomasz Imielinski and Arun N. Swami, "Data Mining a Performance Perspective", IEEE Transactions on Knowledge and Data Engineering, 1993, Vol.5, pp.914-925.
4. Rakesh Agrawal and Ramakrishnan Srikant, "Fast Algorithms for Mining Association Rules in Large Databases", Proceedings of the Twentieth International Conference on Very Large Databases, Santiago, Chile, 1994, pp.487-499.
5. Jiawei Han, Yongjian Fu, "Mining Multiple Level Association Rules in Large Databases", Proceeding of 1995 Int'l Conf. on Very Large Data Bases (VLDB'95), pp. 420-431, Zurich, Switzerland, 1995.
6. R. Agrawal, T. Imilienski, A. Swami, "Mining association rules between sets of items in large databases", Proceedings of the ACM SIGMOD Int'l Conf. on Management of Data, Washington, DC, 1993 pp. 207-216.
7. Han, J., Pei, J., Yin, Y, "Mining frequent Patterns without Candidate Generation" In: ACM-SIGMOD, Dallas, 2000.
8. F. Berzal, J.C. Cubero, N. Marin, J.M. Serrano, TBAR: An Efficient Method for Association Rule Mining in Relational Databases, Elsevier Data and Knowledge, Engineering, pp. 47-64, 2001.
9. D.W. Cheung, J. Han, V.T. Ng, A.W. Fu, Y. Fu, A fast distributed algorithm for mining association rules, Proceedings of Int'l Conf. on PDIS'96, Miami Beach, Florida, USA, 1996.
10. R. Agrawal, R. Srikant, "Mining sequential patterns", Proceedings of the 11th Int'l Conf. on Data Engineering (ICDE), 1995.
11. Ben Franklin, Genealogical Data Mining, 2006.
12. Show-Jane Yen and Arbee L.P.Chen, "A Graph based Approach for Discovering Various Types of Association Rules", IEEE Transaction on Knowledge and Data Engineering, Vol.13, No.5, pp. 839-845, Sep/Oct. 2001.
13. R.Srikant and R.Agrawal, "Mining Generalized Association rules", In Proc. Of the 21st Int'l Conf. on Very Large Databases, pp.407-419, Zurich, Switzerland, 1995.]
14. Cakmak. A, Ozsoyoglu. G, "Taxonomy superimposed graph mining", in proceedings of the 11th Int'l Conf. on Extending Database Technology, Advances in Database Technology, ACM Int'l Conf. Proceeding Series, Vol.261, pp.217-228, 2008.
15. Silvescu. A, Caragea D, Anna Atramentov Graph Databases Artificial Intelligence Research Laboratory, Department of Computer Science Iowa State University Ames, Iowa, 2007.
16. Sadhana Priyadarshini and Debahuti Mishra, "A Hybridized Graph Mining Approach", in proceeding of ICT2010, CCIS 101, pp.356-361, 2010.
17. Haiwei Pan, Xiaolei Tan, Qilong Han, Guisheng Yin, "An Effective Algorithm Based on Association Graph and Matrix for Mining Association Rules", in Proceeding of IEEE2010, 978-1-4244-6977-2/10, 2010.
18. Muthukumar A, Nadarajan R, "Efficient and Scalable Partition based Algorithms for Mining Association Rules", Academic Open Internet Journal ISSN 1311-4360, Vol.19, 2006.
19. Balaji Raja.N and Balakrishnan.G, "Evaluation of Rule Likeness Measures for a Learning Environment", In: Proceedings of the ICOREM Int'l Conf., pp: 1682 - 1690, 2009.
20. Balaji Raja.N and Balakrishnan.G, "A Model of Algorithmic Approach to Itemsets Using Association Rules", In: Proceedings of IEEE, ICECT2011, 978-1-4244-8677-9, pp.6-10, 2011.

