



SURVEY ON DATA MINING APRIORI ALGORITHM

Reeti Trikha

Student, Computer Science Department
Punjab Technical University
India
reeti.trikha31@gmail.com

Jasmeet Singh

Computer Science Department
Punjab Technical University
India
jasmeetgurm@gmail.com

Abstract - Data Mining is known as Knowledge Discovery in Databases. Basically it refers to the extraction of previously unknown and useful information from data in large databases. Association rule mining is a widely used approach in data mining. It helps in revealing interesting relationships. Decision making and understanding the behavior of the customer has become challenging problem for organizations. So for this a technique has been introduced in data mining known as Market Basket Analysis. The algorithm used for learning association rules is Apriori Algorithm. This algorithm finds the frequent pattern based on support and confidence measures. These measures limit the generation of interesting patterns. It's a simple algorithm but having many drawbacks. This paper shows a Survey on working of Apriori algorithm using Tanagra.

Keywords - Data Mining, KDD, Association, Apriori, Market Basket Analysis, Support, Confidence, Tanagra.

I. INTRODUCTION

The main purpose of data mining [4] is to disclose the hidden information from the database. As we know the growth of data volume is increasing day by day in almost every sector like banking, marketing, medicine, website store design, telecommunication, manufacturing, transportation etc.. So data mining proposes a different technique for deletion of repetitive data and conversion of data to more usable form. Basically data mining has two objectives:

PREDICTION: It helps in predicting future values or unknown values of selected variables.

DESCRIPTION: It helps in describing general properties of the existing data in form of patterns.

Data mining is also known as Knowledge Discovery [2] in Databases. Basically they are synonyms for each other. Data Mining is a part of the knowledge discovery process. So, discovery of patterns depends upon the data mining tasks that have been employed. Working of KDD process has been shown in the figure 1.1 below:

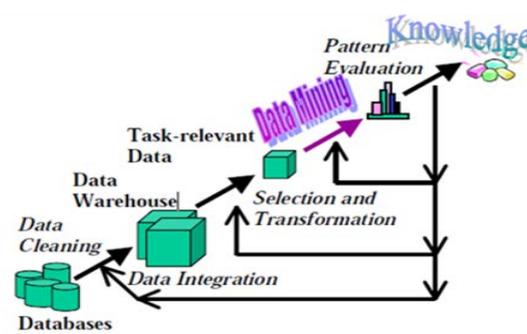


Fig 1.1: KDD Process

Usually it involves three processes as:

First is called PREPROCESSING. This process is executed before data mining techniques are actually applied on the data. It performs data cleansing, data integration, data selection and transformation. Second is called DATA MINING, the main process of KDD process, in which different algorithms are applied to extract hidden information. Then comes third process called as POST PROCESSING, which evaluates the mining result according to requirement of user and domain knowledge. If results are satisfactory, knowledge can be presented, otherwise more processes have to run till we get satisfactory results.

Data mining has discovered many techniques amongst which association rule mining is very important. The main goal of Association Rule Mining [5] is to detect relationships between databases. This is a common task which is helpful in various data mining projects. This introduces one of its best algorithms for mining known as Apriori. The Apriori algorithm [8] discover the frequent patterns from database which support the minimum threshold i.e. it discovers those patterns whose support and confidence must satisfy the minimum support and confidence.

II. ASSOCIATION RULE MINING

2.1 INTRODUCTION

Association rule learning is a method of data mining which helps in discovering interesting relations between variables in a large database. The discovery of frequent patterns is done among a large set of data items. From frequent pattern we



mean those patterns which support minimum support and threshold [3]. For example, a super market analysis indicates that if a customer buys bread and butter together, then there is high probability of buying milk with it. This kind of information helps business organisations to analyse the behaviour of their customers. Shopping centres are using this approach these days to increase their sales by placing the items next to each other. Association Rule Mining has been used in various applications like banking where it helps banks in detecting those customers which are making profit for their organisation. It also helps in detecting frauds to the bank.

2.2 CONCEPTS OF ASSOCIATION RULE MINING

Suppose we have I as a set of items, D as a set of transactions, then association rule is an implication of the $X \Rightarrow Y$, where X, Y are subsets of I , and X, Y do not intersect. Accordingly, each rule comprises of two measures- support and confidence. *Support*: It's the probability that both X and Y will come together in a transaction.

$$\text{Support}(X, Y) = \frac{n(XUY)}{N}$$

N = Total no. of transactions.

Confidence: It's the probability that follows a condition that is a transaction having X also contains Y .

$$\text{Confidence}(X, Y) = \frac{\text{support}(XUY)}{\text{support}(X)}$$

Itemset: It is a set of items in a transactional database.

k-itemset: It is a set of k -items in a transactional database.

Consider the following Transactional database Table-I:

TABLE I
TRANSACTIONAL DATABASE

TransactionId	Pencil	Sharpener	Eraser
1	1	1	0
2	0	0	1
3	0	0	0
4	1	1	1
5	0	1	0

In Table I, 1 represents the presence of item and 0 represents the absence of items. Now let's calculate the following: Consider X =Pencil and Sharpener, Y = Eraser.

It's a 3-itemset where,

$$\text{Support}(X) = \frac{1}{5}$$

$$= 0.2(20\%)$$

$$\text{Confidence}(X \rightarrow Y) = \frac{0.2}{0.4}$$

$$= 0.5(50\%)$$

Support means that pencil, sharpener and eraser all purchased together while confidence says that whenever pencil and sharpener are purchased together, there is also the possibility that eraser is purchased with them.

2.3 MARKET BASKET ANALYSIS

Market Basket Analysis [10] is one of the best examples of association rule mining. It helps in analysing the behaviour of the customer towards the items placed for its purchase by the organisation. Such kind of analysis helps organisations for introducing new strategies which can benefit their organisation. Market Basket analysis assumes that we have a large set of items, e.g. pencil, eraser, sharpener, scale, etc. The probability of buying a particular item or few items or all items will vary from customer to customer. This gives an idea what items people buy together. This strategy is used by marketers to position items and control the way a typical customer traverses their store.

GOALS OF MARKET-BASKET ANALYSIS:

1. It helps in generating association rules.
2. It suggests a way to eliminate duplicate data.

III. APRIORI ALGORITHM

Apriori is one of the association rule mining algorithm which is used to discover all frequent itemsets from transactional database [6]. Apriori is used in many applications. To understand the working of Apriori algorithm we need to be familiar with the following terms:

Transaction: It's an entry in the database which contains a set of items.

Itemset: It's a set of items in a transactional database.

Candidate Itemset (L_i): Items which are only used for the processing initially. They can be all possible combination of itemsets.

Minimum Support: It's a condition which helps to eliminate the non-frequent items from database. It helps in generating frequent items.

Frequent Itemset (Large Itemset (L_i)): The itemsets which contain frequent items are known as frequent itemsets.

Apriori uses iterative approach known as breadth first search (level-wise search), where $k-1$ itemsets are used to generate k itemsets. Apriori works on the following main principle:



“All frequent itemsets should have their frequent subsets”. This means if {c,d,e} is a frequent itemset, then all its subsets {c,d},{c,e},{d,e},{c},{d} and {e} must also be frequent.

Apriori algorithm uses following two concepts:

- (i) Joining and
- (ii) Pruning

Apriori Algorithm Steps:

- (i) Given support threshold ‘s’. Now in first pass, we find the items that satisfy given threshold. The set discovered is called L1, frequent itemset.
- (ii) Now pairs of items in L1, act as candidate pairs C2 for second pass. The pairs in C2 which support s are frequent pairs, L2.
- (iii) The candidate triples, C3 are those sets which have their sets in L2. On the third pass, count the occurrence of triples in C3; those with a count of at least ‘s’ are frequent triples, L3.
- (iv) Proceed as above till we are left with only frequent itemsets.

A general outline of the Apriori algorithm is:

- Firstly, count all itemsets of size K.
- Prune the non-frequent itemsets of size K.
- Generate new candidates of size K+1.
- Again prune the non-frequent itemsets of size K.
- Repeat from step 1, until we are left with no more candidate itemsets,
- When we get all frequent itemsets, generate rules.

Let’s consider an example to show how Apriori algorithm actually works:

Consider a database, D consisting of 9 transactions. Suppose we require min. support count 2. We need to find frequent itemsets first and then association rules will be generated. Table II shows transactional database having 9 transactions.

TABLE II

TID	List of items
T100	11,12,15
T100	12,14
T100	12,13
T100	11,12,14
T100	11,13
T100	12,13
T100	11,13
T100	11,12,13,15
T100	11,12,13

We scan database to identify the number of occurrence of specific item. After that we will compare candidate support count with minimum support count as shown in Fig-3.1 below:

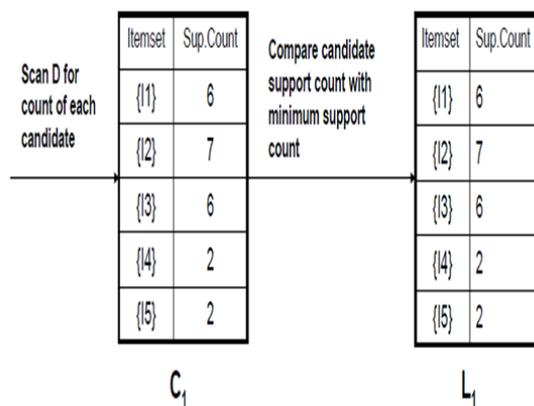


Fig-3.1: Generating 1-frequent itemset

To discover the set of frequent 2-itemsets, L2, the algorithm will use L1. It joins L1 to generate candidate set of 2-itemsets.

Again the transactions in D are scanned and the support counts are compared. The frequent sets of items, L2 are discovered as shown in Fig-3.2

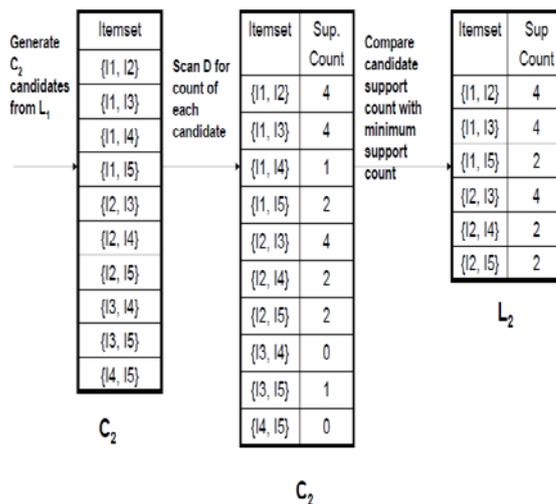


Fig-3.2: Generating 2- frequent itemset

Similarly C3 is computed. We compute L2 and join L2. In this way we are left with frequent itemsets in the end as shown in Fig-3.3



Click on new it will ask for the dataset file as shown in figure-4.2

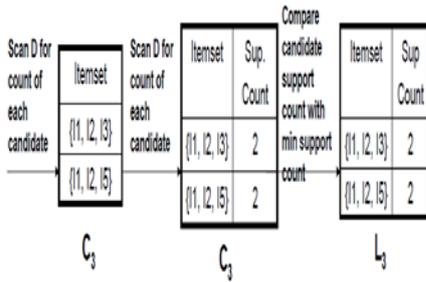


Fig-3.3: Generating 3- frequent itemset

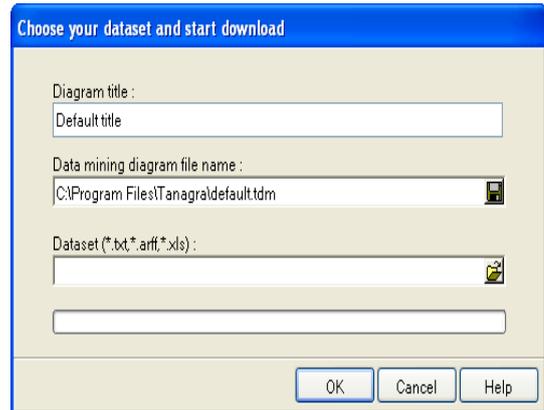


Figure-4.2: Choose Dataset

Choose the dataset by clicking to the dataset icon. It will ask location as in figure-4.3

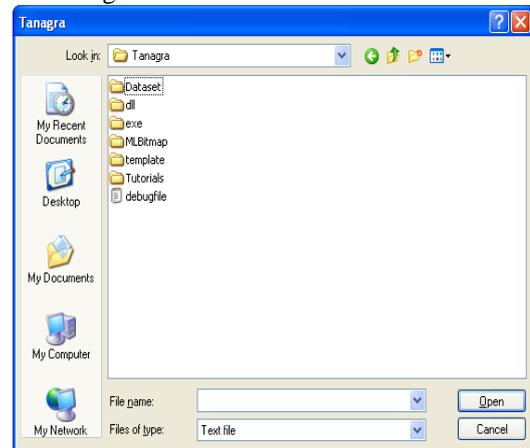


Figure-4.3: Choose Dataset from specified location

After choosing dataset from the folder the path is shown as in figure-4.4

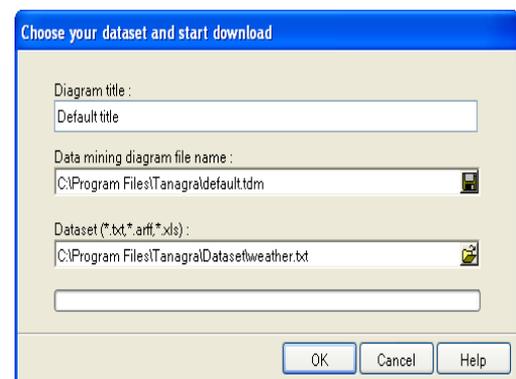


Figure-4.4: Dataset location

Click ok and it will show the Dataset Description as shown in figure-4.5

Pseudo Code For Apriori Algorithm [1]:

C_k :Set of Candidate itemsets of size k
 L_k : Set of Frequent itemsets of size k
 $L_1 = \{ \text{Frequent items} \};$
 For ($k=1; L_k \neq \text{null}; k++$) do begin
 C_{k+1} = candidates generated from L_k ;
 For each transaction in database do
 Increment the count of all candidates in C_{k+1} that are contained in t
 L_{k+1} = Candidates in C_{k+1} with minimum support
 End

Return L_k ;

Apriori algorithm is the classical and simplest algorithm to implement the concept of association rule mining. But there are some disadvantages as follows:

Drawbacks of Apriori Algorithm:

1. It consumes lot of time for scanning database.
2. It increases space complexity by generating a large number of non-frequent itemsets.
3. It makes several iterations while mining data.
4. Generation of large amount of frequent itemset, also makes it inefficient to use.
5. It does not consider the significance of item to the user or business. Its focus is only on the presence and absence of the item.

IV. WORKING OF APRIORI USING TANAGRA

Figure-4.1 will be the first screen when Tanagra starts up

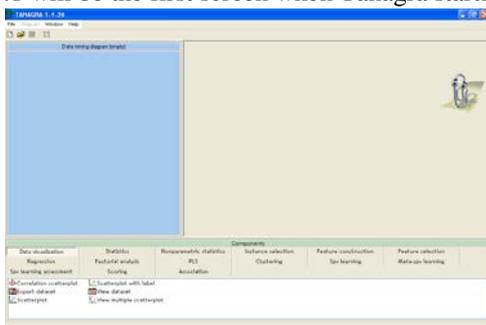


Figure-4.1: Tanagra Home Screen

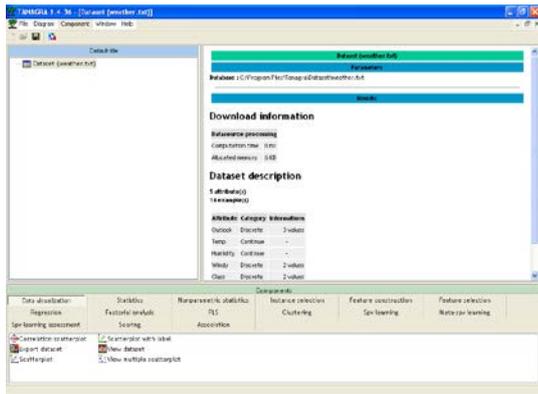


Figure-4.5: Dataset Description

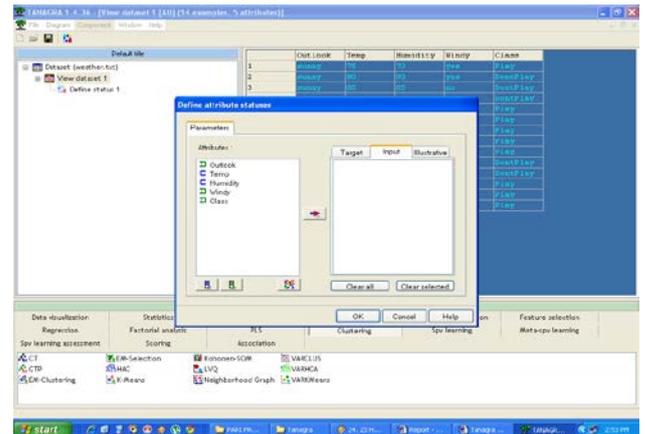


Figure-4.8: Define Status

Drag view dataset from lower window to upper left window onto Dataset as shown in figure-4.6

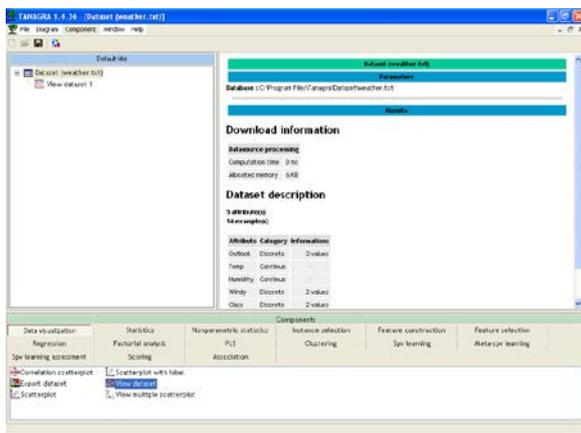


Figure-4.6: View Dataset

Right click on view dataset 1 and click on view, following figure-4.7 appears

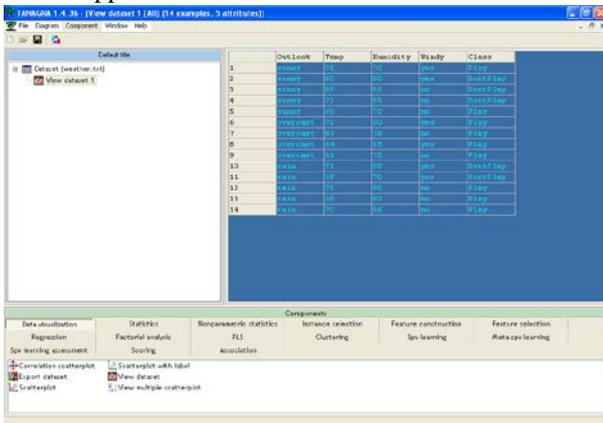


Figure-4.7: Dataset viewing screen

Add define status from toolbar just close to save button as shown in figure-4.8

For Apriori algorithm we need two or more discrete attributes as input and no need for target. So select as shown in figure-4.9. For input and target attributes we can have help in form of tooltip by placing mouse over the respective algorithm.

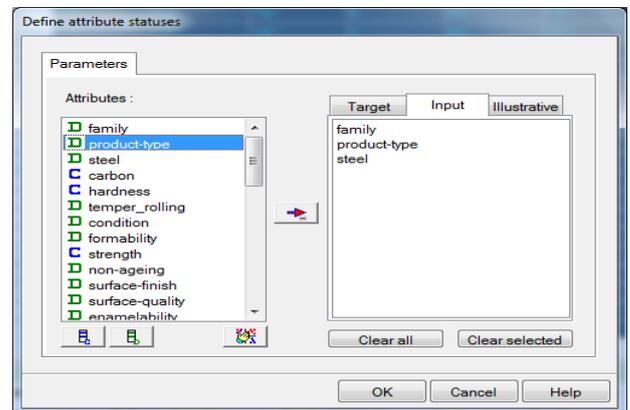


Figure-4.9: Define Status attributes

Click on association from lower window drag Apriori over define status as in figure-4.10

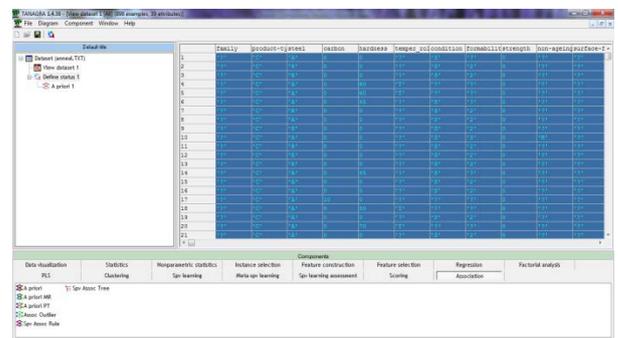


Figure-4.10: Apriori

Right click on Apriori, execute and view the result as in figure-4.11

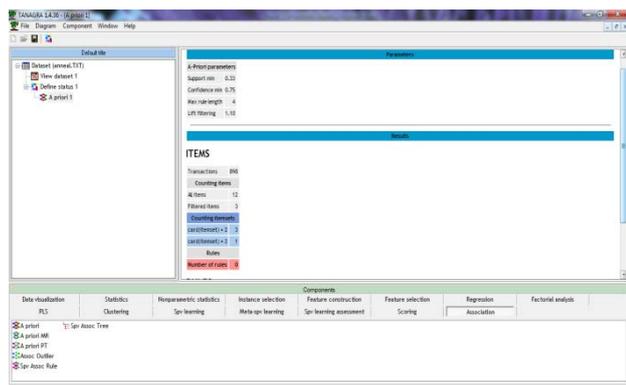


Figure-4.11: Apriori Result

V. CONCLUSION AND FUTURE WORK

Association rule mining is useful in discovering useful patterns from transactional database. This algorithm is used for the implementation of association rule mining. But this classical algorithm has several limitations like scanning time, memory optimization, candidate generation, etc.

Our classical Apriori treats all the items in database equally by focusing only on the presence and absence of an item within the transaction. It ignores the significance of that item to user or business. So, Apriori algorithm efficiency can be improved further to provide valuable information to customer as well as business.

REFERENCES

- [1] Ms. Arti Rathod, Mr. Ajaysingh Dhabariya & Mr. Chintan Thacker, (Sep. 2013), "A review on Association Rule Mining and Improved Apriori Algorithms", International Journal of Scientific Research in Computer Science, vol. 1, no. 11.
- [2] Frawley, W., Piatetsky-Shapiro, G., Matheus, (1992), "Knowledge Discovery in Databases: An Overview", AI Magazine, fall 1992, pp. 213-228.
- [3] Chun-Jung Chu, Vincent S. Tseng and Tyne Liang (November (2008), "Mining temporal rare utility Itemsets in large databases using relative utility thresholds", International Journal of Innovative Computing, Information and Control, Vol. 4, no. 11.
- [4] Jiawei Han and Micheline Kamber-book second edition, "Data Mining Concepts and Techniques".
- [5] Irena Tudor, Universitatea Petrol-Gaze din Ploiesti, "Association Rule Mining as a Data Mining Technique", Bd. Bucuresti 39, Ploiesti, Catedra de Informatica, Vol-LX, No.1, 2008.
- [6] Mamta Dhanda, (July 2011), "An efficient approach to extract frequent patterns from transactional database", International Journal of Engineering Science & Technology, vol. 3, no. 7.
- [7] Wei Zhang, Hongzhi Liao and Na Zhao (2008), "Research On The Frequent Pattern Growth Algorithm about Association Rule Mining", International Seminar on Business and Information Management.

[8] Zhang Shu-mao and DU Ying-shuang (Sep 2008), "The analysis and improvement of Apriori Algorithm", HAN Feng Journal of Communication and Computer, vol. 5, no. 9.

[9] Du Ping and Gao Yongping (2010), "A New Improvement of Apriori Algorithm For Mining Association Rules", International Conference on Computer Application and System Modeling.

[10] Dr. Dhanabhakya, Dr. M. Punithavalli, Dr. SNS college of Arts and Science, "The Survey on Data Mining Algorithm for market basket analysis", Global Journal of Computer Science and Technology (IJCSIT), Vol. 11, Issue 11 Version 1.0, July 2012, ISSN: 0975-4172.

[11] Sheila A. Abaya, "Association Rule Mining based on Apriori algorithm in minimizing candidate generation", International Journal of Scientific & Engineering Research, Vol. 1-3, Issue-7, July 2012, ISSN 2229-5518.

[12] Raorane A.A, Kulkarni R.V and Jitkar B.D "Association Rule- Extracting Knowledge Using Market Basket Analysis", Dept. of Computer Science, Vivekanand College, Tarabai Park Kolhapur, Research Journal of Recent Sciences, Vol. 11(2) 19-27, Feb. 2012.