

# Augmentation of Apriori Algorithm

<sup>[1]</sup>Nivedita, <sup>[2]</sup>Jhelum Dasgupta, <sup>[3]</sup>Razauddin  
<sup>[1][2][3]</sup>Master's Scholar, Dept. of Computer Sc. & Engineering  
Birla Institute of Technology  
Mesra, India

<sup>[1]</sup>nivepriya09@gmail.com, <sup>[2]</sup>Jhelum.dasgupta@gmail.com, <sup>[3]</sup>razauddin99@gmail.com

**Abstract:** There is hasty evolution of system and information technology, ceaseless information has gained people attention progressively. While there is need of information at a fast paced flow, the investigation and mining of the information along with the regulations deeply jiggled among statistics are of prime importance. Data mining knowledge is to formulate and consider the information, which extracts and find awareness since the piles of statistics, therefore how to define the information mined.

Apriori algorithm is in trend and a typical method in data mining. The basic proposal of assumed procedure is in the direction to have a helpful model in wide-range of statistics sets. The procedure lacks in numerous domains. This research deals through the apriori algorithm, and varied procedure is being projected to enhance the strength of apriori algorithm. The algorithm is basically used for association rule mining field.

**Keywords**— Apriori algorithm, Association Rule mining, Data mining

## I. INTRODUCTION

Modern computer & database technology only clusters enormous information and is not able to adequately assemble and utilize the acquaintance lying within them. The evolution of data mining tools met these people needs. The piles of data specifically in network service-based applications, service-leaning design, and cloud computing, field is multiplying consistently and to find out the valuable information from it there is need of various data mining technologies. Data mining discern useful information from the data. Association rule mining, carve out interesting relationships or correlation association amid certain objects in a huge database. It is an essential technique in the information mining. Association study is valuable for finding out motivating relationships unseen in huge records sets. The discovered associations can be represented in the variety of relationships policy or sets of repeated objects. Association shows the dependencies or the associations between various data sets available. It tries to locate and identify the embedded or interesting association in the statistics set, and the outcome is generally put up in the shape of common data sets or association rules. The matter of worry in recent time about Relationships rule mining is to develop the algorithmic potential and performance. Association law mining is to find out the probable relation among sets of records objects.

Apriori algorithm is an individual of basic procedure in relationships rule mining. Apriori algorithm was initiated via 'Agarwal' in 1993. This procedure is one which best

describes the association mining. Apriori is one of the classical methods used for relationships rule generation. Its main idea focuses on identified low-dimensional common item sets to figure out higher dimensional repeated item sets. A unique and essential asset of apriori algorithm is, "of an item set is not common, and some of its superset is not at all common". This feature of apriori discovers frequent item sets. With the evolution of Internet, data amount has increased many folds; the conventional Apriori clustering algorithm is not in accordance to meet the ever changing demands. Apriori, uses a breadth first search (BFS) technique. Apriori is an influential algorithm engaging a continual accession identified as a level-wise and BFS, which k-data sets are used to create (k+1)-data sets. The fundamental and still superior, Apriori feature called anti-monotone which having the property of efficiently generating candidate item sets by taking out unnecessary items. Apriori has two-pace process, join as well as prune. Some of the common terminologies associated with Apriori algorithm are:

**Item set** - Item set is compilation of items in a database being denoted by

$$I = \{i_1, i_2, \dots, i_n\}, \text{ where } n \rightarrow \text{number of items.}$$

**Transaction** - Transaction is a database entry having assemblage of items. Transaction  $T$  and it is  $T \subseteq I$ . A transaction having item set

$$T = \{i_1, i_2, \dots, i_n\}.$$

**Minimum support** - Minimum support condition need to be fulfilled by the given items such that added processing of

items being carried out. Minimum support could be esteemed as a happening condition that removes the in-frequent items of any database. Customarily the bare minimum support is in percentage form.

**Frequent data set** – The item sets which having the bare minimum support criteria are to be referred as frequent data sets. Denoted by  $L_i, i \rightarrow i^{\text{th}}$  data set.

**Candidate item set** – Candidate item set are those items that are only considered for processing. Candidate item set is identifying every possible combination of item set. Denoted by  $C_i, i \rightarrow i$ -item set.

**Support** –Support threshold represents the usefulness of a rule. Support measures number of transactions having item sets that match either sides of the implication in the association rule. Let's figure out, two items A and B and calculate support of  $(A \rightarrow B)$  formula to be used is: -  
Support  $(A \rightarrow B) = (\text{amount of transaction enclosed by both A as well as B}) / (\text{Total amount of transaction})$ .

**Confidence** –It demonstrates the inevitability of the rule. This parameter counts how frequently a transaction's item set maps in accordance with its left side implication in the association rule matching it for right side. The item set not satisfying above both situations can be useless. Let's take two objects X as well as Y and calculate assurance of  $(X \rightarrow Y)$ , subsequent formula is used: -

Conf  $(X \rightarrow Y) = (\text{number of transaction enclosed by both X as well as Y}) / (\text{Transaction enclosed only by X})$ .

**Apriori Algorithm** - Apriori method is based on 2 concepts functioning namely:-

a) Self Join and b) Pruning. Apriori utilizes level wise exploring,

$X \rightarrow$  item sets worned to discover  $(X+1) \rightarrow$  object sets.

1. The set of frequent 1- item sets being found and is depicted as  $C_1$ .

2. The next step is calculation of support that means the occurrence of the item in the database. This scans the whole database.

3. Then pruning is performed on  $C_1$  in which items mapped with the minimum support parameter. The items satisfying the minimum support criterion are just considered for further processes that are denoted by  $L_1$ .

4. Further candidate set generation step which generates 2-itemset which are denoted by  $C_2$ .

5. Further database is browsed for calculation of 2-itemset support. As per the minimum support the generated candidate sets are examined and item set satisfying the minimum

support criteria are carried out further for 3-itemset candidate set generation.

6. These above steps continue by the time no frequent or candidate set is generated.

## II. RELATED WORK

The responsibility of discovering relation or u can say connection and patterns emerged when the current market sense the urgency to learn apprentice user demeanors of procuring, for sales boost up, raise business, engage more customer attention, and to boost the profit. So, in 1993, R. Agrawal et.al sense urgency to design and develop a advance and more efficient algorithm which is known as Apriori Algorithm, for identifying and finding relation among the dataset. Fast algorithms for mining association rules' by R. Agrawal and R. Srikant generated in 1994. In this execution time improved considerably as dataset number increases. [1] In the year 1995, R. Agrawal came up with an advance algorithm for mining association rules for large dataset. This algorithm reduced CPU aloft and in this Partitioning Algorithm is being defined. [2] In 'Mining Sequential Patterns', the complete operation of mining was categorized into five phases: sort phase, L itemset stage, transformation phase, sequence phase and maximal phase. Apriori algorithm is a classical algorithm in the association rule mining area; the basic idea is to make use of a low-dimensional frequent item sets to reap in high-dimensional frequent item sets by continuous iterative steps. Various scholars have remodelled and upgraded the Apriori algorithm by different view as of now. Brin contemplated a dynamic item sets counting algorithm DIC [3], which reduces the frequency of scanning the database very effectively. This algorithm divides the transaction database of same size, and accesses data blocks consecutively to generate the 1-frequent item sets then generating the candidate 2-frequent item next thus sets by self-join, lastly merging the 1-frequent item sets and candidate 2-frequent item sets that has been generated by each block, now repeating it till no new item sets or has reached the limit. The DHP algorithm contemplated by Park et al. [4] which has idea about the dynamic hash hashing algorithms and pruning algorithm, and it excludes transactions not generating frequent item sets when scanning the database, thus improving the mining efficiency of frequent item sets. For big data or where the aspect of data is too high. Google turned its focus towards Map Reduce [5] framework in 2004, then Hadoop which depend on Map Reduce became mainstream thing, and cluster-based parallel data mining garnered much of the focus, research and utilization. Hadoop is used for parallelizing many classical data mining algorithms. Stationed around Apriori algorithm, Lin et al suggested various parallel algorithms namely - SPC,

FPC and DPC [6] based on Map Reduce, the SPC algorithm allocate the data set to all Map nodes, thus executing mining operation parallel, next the reduce phase performing joining operation, the algorithm begins with the Map and Reduce work once, nonetheless the FPC and DPC algorithms pressing to begin Map Reduce tasks repetitively, as decisive by the dimension of frequent item sets mining; Li et al contemplated alongside recurrent item sets mining algorithm PApriori [7], in Map stage going in through the transaction database to calculate candidate repeated entity sets, and execute statistical operations to have repeated item sets in reduce phase, despite that it also need to repeatedly start Map Reduce tasks.

### III. CLASSICAL APRIORI ALGORITHM

Implementing an insistent approach, in every insist Apriori algorithm develop candidate item-sets using large itemsets of a previous insist. Key perception of this iterative advent is as listed below:

#### Apriori\_algorithm( $L_i$ )

1.  $K_1 = \{\text{frequent-1 item-sets}\}$ ;
2. for ( $i=2; K_{i-1} \neq \Phi; i++$ ) {
3.  $D_i = \text{discover\_Apriori}(L_{i-1});$  //novel candidates
4. for all transactions  $x \in C$  do begin
5.  $D_i = \text{subset}(D_i, x);$  //Candidates
6. for all candidates  $c \in D_x$  do
7.  $c.\text{count}++;$  }
8.  $K_i = \{c \in D_i \mid c.\text{count} \geq \text{minsup}\}$
9. end for
10. Answer= $G_i K_i$

### IV. PROPOSED ENHANCEMENT IN EXISTING APRIORI ALGORITHM

#### 1. Improvement of Apriori

In this dawn to augment Apriori algorithm efficiency, we focus on shortening the time expend for  $C_k$  generation. In the action to locate frequent item sets, first size of a transaction (ST) is raised for every transaction in DB and maintained. Now, locate  $L_1$  containing set of items, support value for each individual item and transaction id's having the item. Use  $L_1$  to generate  $L_2, L_3 \dots$  clubbed with diminishing the database size such that instance time is reduced to examine the transaction since the database. To create  $C_2(x, y)$  (objects in  $C_k$  be  $x$  and  $y$ ), do  $L_{(k-1)} \times L_{(k-1)}$ . To discover  $L_2$  as of  $C_2$ , in preference of examining entire record and each

and every one of its communication, we get rid of operations where  $ST < k$  (where  $k$  is 2, 3...) and also get rid of the removed operations from  $L_1$  the same as well. This facilitates in diminishing the moment to examine through the uncommon transactions commencing the database.

Discover the bare minimum support commencing  $x$  as well as  $y$  and obtain transaction id's of bare minimum support count entity commencing  $L_1$ . At the moment,  $C_k$  is browsed in favor of précised transactions simply and with reduced DB magnitude. After that,  $L_2$  brings out by  $C_2$ , where support of  $C_k \geq \text{min\_support}(x, y, z)$ ,  $L_3$  and so forth is brought out by to reproduce the above process until and unless no frequent objects sets can be revealed.

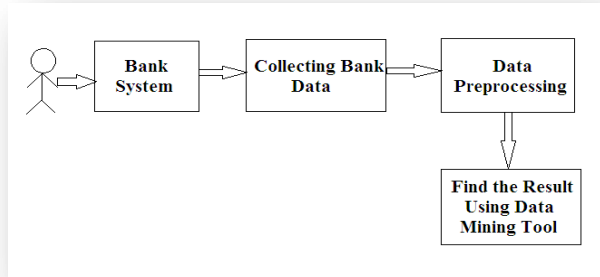
#### Apriori Algorithm

**Input:** transactions record,  $N$

Minimum support,  $\text{min\_supp}$

**Output Lk:** frequent entitysets in  $N$

1. Get ST //for every transaction in DB
2.  $L_1 = \text{find\_frequent\_1\_entityset}(N)$
3.  $L_1 += \text{get\_txn\_ids}(N)$
4. for ( $k=2; L_{k-1} \neq \Phi; k++$ ) {
5.  $C_k = \text{generate\_candidate}(L_{k-1})$
6.  $x = \text{item\_min\_supp}(C_k, L_1)$  //find object from  $C_k(i, j)$  which have minimum support by means of  $L_1$
7.  $\text{target} = \text{get\_tx\_id}(x)$  //find transactions for every item
8. for each ( $tx\ t$  in  $tg$ ) do {
9.  $C_k.\text{count}++$
10.  $L_k = (\text{entity in } C_k \geq \text{min\_sup})$
11. } //finish for each
12. for each ( $tx$  in  $N$ ) {
13. if ( $ST = (k-1)$ )
14.  $tx\_set += tx$
15. } //finish for each
16.  $\text{delete\_tx\_DB}(tx\_set)$  //decrease DB size
17.  $\text{delete\_tx\_L1}(tx\_set, L_1)$  //decrease transaction size in  $L_1$
18. } //finish for



**Figure - 1: Framework to find the Association Rules**

**Apriori Pseudo code**

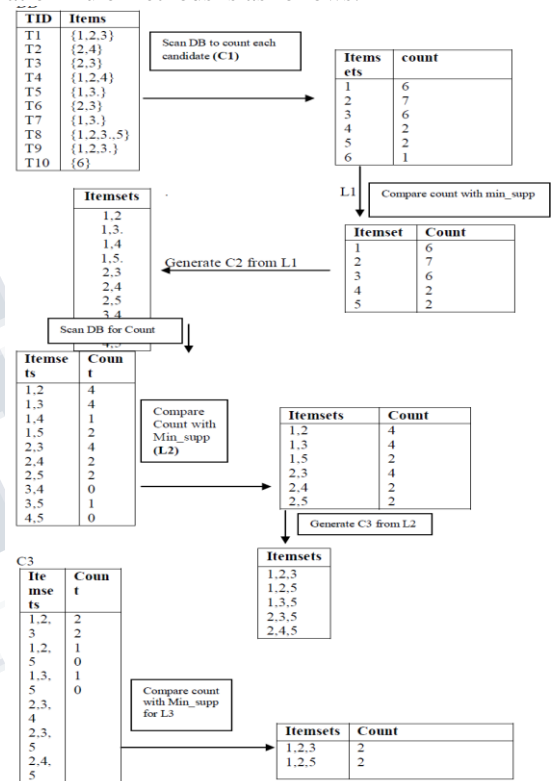
1. Apriori (T,ε)
2.  $L_1 \leftarrow \{ \text{large 1-itemsets that appear in more than } \epsilon \text{ transactions} \}$
3.  $j \leftarrow 2$
4. while  $L_{j-1} \neq \emptyset$
5.  $C_j \leftarrow \text{Generate } (L_{j-1})$
6. for transactions  $t \in T$
7.  $C_t \leftarrow \text{Subset } (C_j, t)$
8. for candidate's  $c \in C_t$
9.  $\text{count}[c] = \text{count}[c] + 1$
10.  $L_j \leftarrow \{ c \in C_j \mid \text{count}[c] \geq \epsilon \}$
11.  $j \leftarrow j+1$
12.  $UL_j$
13. return  $j$

**2. Essential Terms Utilized in Apriori**

- a) Minimum\_support: it is bare minimum support utilized for penetrating common patterns that can gratify this limitation.
- b) Minimum\_confidence: it is bare minimum confidence utilized for discovering the powerful association imperative that can gratify this threshold
- c) Frequent\_Itemset ( $L_k$ ): represented by  $L_k$ , where  $k$  means  $k^{\text{th}}$  item, these are the entity sets that gratify the bare minimum support (min\_support) threshold.
- d) Join\_Operation: on behalf of  $L_i$ , a set of candidate i-itemsets ( $C_i$ ) is created by the union of  $L_{i-1}$  with  $L_{i-1} (L_{i-1} \times L_{i-1})$
- e) Apriori\_Property: practical assets for frill inappropriate data. It specifies; some subset of frequent entityset should be recurrent.
- f) Prune\_step: utilized in for discovering frequent datasets, for any (k-1)-datasets that is not common can't turn out to be subset of a common i-itemset.
- g) Definitions:  $L_i$  – set of common itemsets of "i" size created by utilizing min support.  $C_i$ – set of candidate entitysets of "i" size.

**3. Relative Learning of Association Rule Methods :**

We have taken into account 3 association rule procedure i.e. Predictive Apriori association rule, Apriori Association rule as well as Tertius association rule. And compared there outcome of those association rule methods with the aid of information mining tool. Apriori Association Rule method is elucidated in the preceding subsections. At the moment, short explanation of extra two algorithms i.e. Predictive Apriori Association Rule methods and Tertius Association Rule methods is as follows:



**Figure - 2 : Creation of Candidate entitysets and Frequent Itemsets (Min\_support=2 (20%))**

**Predictive apriori association rule method:** apriori association rule algorithm, this analytical accurateness is warned to produce the apriori association rule. In weka, the method generates "q" most excellent association rule which stands on "q" specified by the user.

**Tertius association rule algorithm:** the method looks to discover the rule according to the authentication actions. It makes utilization of 1st order logic demonstration. It comprises of diverse choices or parameters like frequency threshold, values, class index, horn clauses, classification, confirmation values, confirmation threshold, missing



negation, noise threshold, roc analysis, number literals, Values Output, Repeat Literals, Etc.

### V. EXPERIMENTAL RESULT

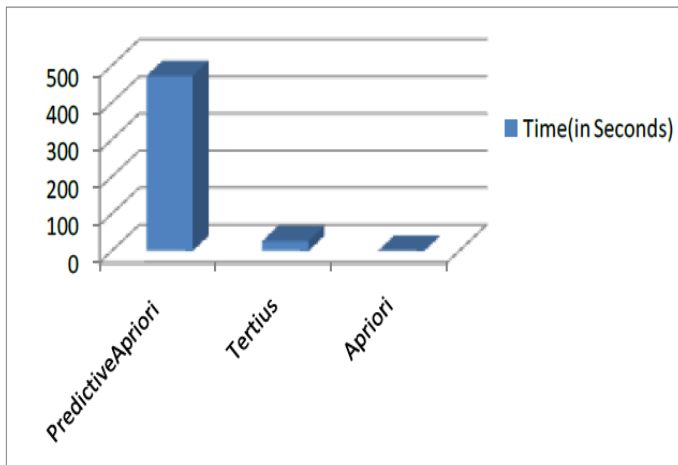
We got the results by the use of these three relationships rule methods. In the paper, we have used bank information for assessment by 11 attributes as well as 600 data and found association rules utilizing WEKA. Table 1 signifies the outcome by utilizing Predictive Apriori Relationships Rule Method. 1st column signifies important attributes record later than filtering the unimportant attributes in the illustration customer id is deleted. Table 2 signifies the outcome by utilizing Tertius Association Rule Method. Table 3 signifies the outcome by utilizing Apriori Association Rule Method. The assortment of income attribute on behalf of amenity is described as follows: -

```
"(-inc-24386]"→'(-ing-24386.173333]',
"(24386- 43758]"→ "(24386.173333-43758.136667]",
"(43758-inc)"→'(43758.136667-ing)'
```

And age attribute is represented as:

```
"(-inc-34]"→'(-ing-34.333333]',
"(34-50]"→'(34.333333-50.666667]',
"(50- inc)"→"(50.666667-ing)"
```

Below given Figure- 3: gives a comparison time taken by 3 methods:



**Figure - 3: Comparison of Predictive Apriori, Apriori and Tertius on the foundation of Elapsed time**

### VI. CONCLUSION

We boosted the idea of the Apriori method efficiency by minimizing the time used to examine through the DB transactions. We found with the aim of that k value being increased, number of transactions inspected reduces and therefore, time consumed in addition takes a dip in analogy to traditional Apriori method.

Due of this, instance taken by spawn candidate data sets with our suggestion also takes a dip as compared to classical Apriori.

### FUTURE SCOPE

Discussed methodology can be utilized in further domains to carry forward the significance among the statistics there in the depository. Association rules created by these three methodologies could be merged together for improved consequences of some valid living application. Algorithms might also be pooled to form an capable algorithm.

### REFERENCES

- [1] R. Agrawal, and R. Srikant, “Fast Algorithms for Mining Association Rules”, In Proc. VLDB 1994, pp.487-499
- [2] A. Savasere, E. Omiecinski, and S. Navathe, “An Efficient Algorithm for Mining Association Rules in Large Databases”, In VLDB’95, pp.432-443, Zurich, Switzerland
- [3] S. Brin, et al. Dynamic item set counting and implication rules for market basket data. Proceedings of the ACM SIGMOD International Conference on Management of Data. 1997. 123-140
- [4] J. S. Park, M. S. Chen, P. S. Yu. Efficient parallel data mining of association rules. 4th International Conference on Information and Knowledge Management, 1995, 11: 233-235P
- [5] Jeffrey Dean, Sanjay Ghemawat. Map/Reduce: Simplified Data Processing on Large Clusters[R]. OSDI’04: Sixth Symposium on Operating System Design and Implementation 2004.
- [6] Li N., Zeng L., He Q. & Shi Z. (2012). Parallel Implementation of Apriori Algorithm Based on MapReduce. In: Proceedings of the 13th ACM International Conference on Software Engineering, Artificial Intelligence, Networking

and Parallel & Distributed Computing (SNPD '12). Kyoto, IEEE: 236–241.

[7] Lin M., Lee P. & Hsueh S. (2012). Apriori-based Frequent Item set Mining Algorithms on MapReduce. Proc. of the 16th International Conference on Ubiquitous Information Management and Communication (ICUIMC '12). New York, NY, USA, ACM: Article No.76.

