

Mining Frequent Itemset on Temporal Data by Using Improved Apriori Algorithm

¹Gayatri N. Rukare

¹Department of Computer Engineering
Gokhale Education Society's R. H. Sapat college of Engineering,
Management Studies and Research, Nashik-05

Abstract: Now a days, Data mining is an important research area. But the frequent itemset mining is a part of data mining in which there are extensively improving. There are different techniques available for getting frequent itemset by using the different rules. Apriori and FP growth algorithm also works on the frequent itemset mining. Apriori takes too much time to produce the output so that it's efficiency is less. FP growth algorithm relayed on the searching, sorting and many more. In FP Growth algorithm header table is used to store the transaction ID and transaction of the dataset. Header Table plays vital role in creating new data structure which leads to create a master table which is called as enhanced header table which stores the frequent transaction with the transaction Id to improve the efficiency of mining the frequent itemset.

Index Terms: Data mining, frequent itemset, header table , header tree, apriori algorithm, Fp-Growth algorithm, temporal data, frequent pattern.

I. INTRODUCTION

Data mining is the important research area. In data mining the frequent item mining is the key factor in which all the information is compressed. After that it will scan the data and form the master table with the transaction ID. Each time it will scan the all the data set for all the items in standard apriori algorithm. In FP-Growth Algorithm all the item Scanned and form the FP Tree with the different level using the binary search tree. Now in proposed system all the data is scanned once and formed the header table with the transaction ID. When next iteration is performed it will not scanned all the database it will scanned the only new data item if it finds the new item it will add it in to the header table. This process reduced the time required for the finding frequent itemset then generate the candidate key for the comparison with the large 1 item in each iteration new entry is added to the header table. Then all the entity of the header table form the header tree with different level e.g level 1, level 2.... etc. then it will form n-pattern. Then it will calculate the frequent itemset for sorting the itemset it will use the binary search tree algorithm.

An efficient algorithm to mine frequent patterns and their related time interval from transactional database. The header tree is developed for reducing the time required for mining the frequent itemset according to the dataset. Then a new algorithm is proposed based on two thresholds, support, and density as a novel threshold. Frequent itemsets are discovered and merged with the those with neighboring time intervals.

II. RELATED WORK

Rakesh Agrawal et al. [2] They presented an algorithm which generates all association rules between frequent items in the data set. The algorithm cannot handle the buffer management and pruning techniques. The algorithm had best performance. The estimation procedure exhibited high accuracy and the pruning techniques were able to prune out a very large fraction of itemsets without measuring them.

R. Srikant et al. [3] They present two algorithms, Cumulate and EstMerge, which is having high efficiency than Basic association rule mining. They presented a new algorithm to find the frequent items along with the association rule mining. It uses the pruning technique to find frequent itemset from the dataset.

Bing Liu et al. [4], They propose two mining techniques. The integration is focused on the Class association rule which is very efficient because of the building a classifier based on set.

Juan M. Ale et al. [5] They have introduced the problem of association rules discovery, given place to what we call Temporal Association Rules. They have also introduced the concept of temporal data support were not discovered with the traditional mining frequent itemset.

Yingjiu Li et al. [6] They defines two types of temporal association rules: precise-match association rules and fuzzy-match. They also worked on the Apriori algorithm, and the calendar-based patterns which is the advance technique.

Yingjiu Li et al. [7], They studied the association rules mining with their temporal patterns in terms of calendar schemas. They studied the Apriori algorithm and extended it with time stamp. They developed two optimization techniques to improve the performance of the data mining process.

Xindong Wu et al. [8] presented an algorithm which works on the both positive and negative association rules. The method extends traditional associations to include association rules of temporal data for large scale database while mining the frequent itemset with respect to time.

Yen-Liang Chen [9] propose a new algorithm to work on weakness Market basket analysis. All the dataset is considered as shelf of the store and stores are having the different size. It will consider the all the data with the time stamping information. On the large number of stores and period is considered for the execution.

Yen-Liang Chen et al. [10] propose a representation scheme to which is not depended on the price of transaction data. An efficient algorithm is developed in two parts in which first is patterns with price information and child pattern can be classified in to different categories.,some of them are positive , negative and neutral values.

Yen-Liang Chen et al. [11], introduce the data mining problem and define the rule patterns that can be mined from different data. A unified approach is developed based on fuzzy techniques so that all different data types can be handled in a uniform manner. After that, an algorithm is developed to discover fuzzy association rules from the questionnaire dataset.

III. PROPOSED METHODOLOGY

A. Apriori Algorithm

Apriori algorithm is well known algorithm for finding the frequent itemset from large data. Frequency count of the each item increased as it scans the dataset. It scans whole data set at every time to find the frequent itemset. It elaborates in the 3 steps such as:

1. It will partition the data in the different itemset with the specific time stamp.
2. Generate candidate key and compare it with the large1 itemset
3. It will generate frequent itemset

B. FP- Growth Algorithm

FP-Growth algorithm uses the FP-tree formula for finding the frequent itemset from large dataset. It compresses the information and associate it to form FP instances to represent the frequent items in the dataset. then it compressed information into a collection of conditional databases, each linked to a common template. Finally , each information is extracted one by one.

C. Improved Algorithm

FP-growth algorithm takes more time to scans the database so here it can use the dynamic binary search tree data structure Factors that affect the performance of the FP-tree algorithm mining process. Inserting items into the tree. Searching for entries, walking through fp structures, data structures, node structures and pruning are few. only The mining process starts with a transactional database scan to find a common set of data sets during this process. data structure Used to store all items will affect the working time of the mining process. The total number of transactions in the transaction database may not be known in the preceding step. To use a simple data structure array, we should know the number of previous entries since the array is a static data structure. To calculate the number of items, we need to scan the database again. The growth of FP and other related methods uses two scans to create a tree.

3.2 System Architecture:

FP tree can be a solid information design that has retained essential, important and quantitative information taking into account common patterns are as follows:

- Contains 1 root marked "root", a group of prefix sub tree, part of the foundation kid, and a frequent header table.
- A one-to-one node in the part prefix sub tree contains 3 fields:
- Item-name: Lists the element representing this node.
- Number: records the number of transactions represented by the portion of the track that returns to the current node.
- Node-link: Connects to the next node in the tree with the same element name, or null if none.

After getting the support count of each item ,a header table is created by using the frequent items only. This time we have a prior knowledge about the number of items. Therefore we can create array here. But to use a static data structure the availability of continuous free memory is an issue. Here again a binary search tree is used as a header tree instead of header table called Binary Search Header Tree(BSHT) . Each node of the BSHT includes item name , support count ,link to the next node and link to the corresponding item in the fptree. If we are using BST while counting the item frequency we can prune the infrequent items and can reuse the same BST as BSHT with a slight modification. The BST is created earlier by considering item names as keys. In BSHT also the item names are used as keys. The structure of a binary search tree depends on the order of items . In this implementation the number of items is available. So to get the maximum efficiency here we created the BSHT after sorting the items in frequency descending order to get the most frequent item in the root of BSHT.

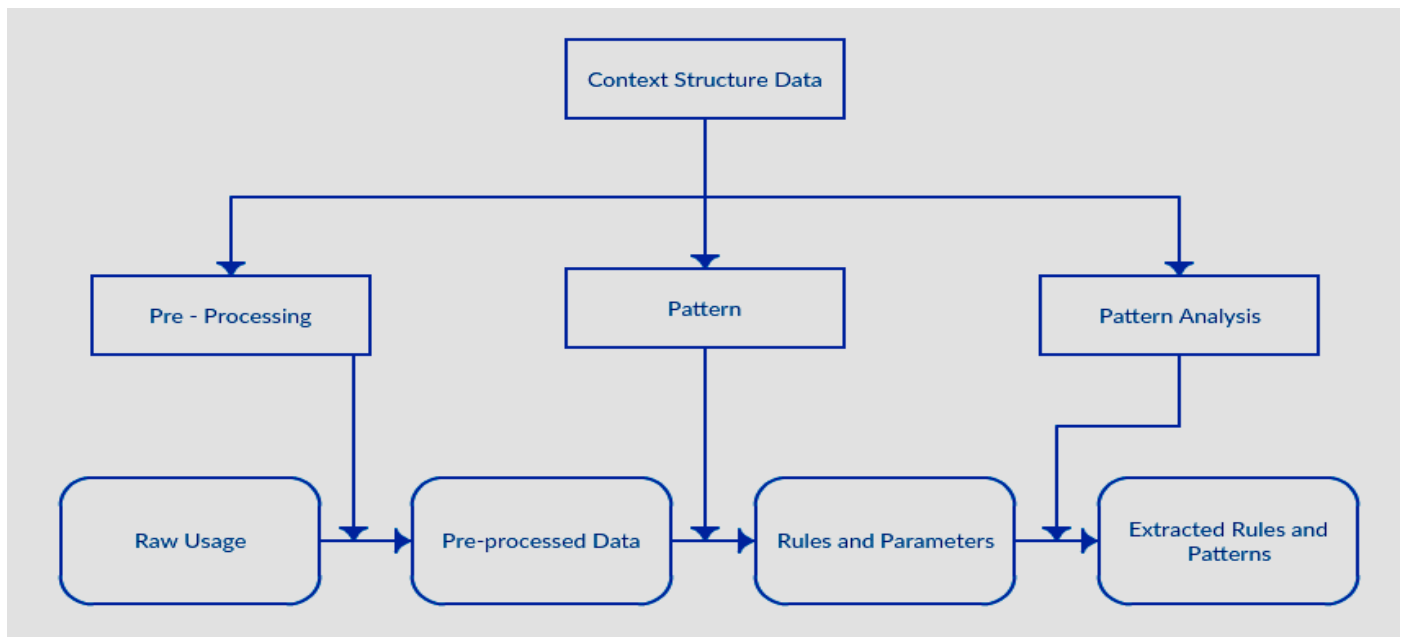


Fig.1 System Architecture

Then the second scan starts for constructing the tree by inserting each transaction one by one. Before inserting a transaction in to the tree the infrequent items have to be removed and the remaining items should be sorted in descending order of the frequency . In earlier algorithms to remove the infrequent items we should search the header table to verify whether the item is there or not. The header table is the set of items sorted in descending order of their frequency.

Each node contains item name , frequency and link to the next node. During the second scan we can reuse the same BST as the header tree, but we have to do two things. First, include another field in the nodes as link to point to the corresponding first node of the item in the fptree, second, delete all infrequent nodes from the BST by using the BST node removal process. We can use array or BST to use as a header table. Both of them perform efficiently and we can apply a binary search on both the data structure. Illustrates the header tree with fp tree for table (a).

| Items | Frequency | Links |
|-------|-----------|-------|
| A | 2 | * |
| B | 4 | * |
| C | 6 | * |
| D | 4 | * |
| E | 3 | * |
| F | 5 | * |
| G | 2 | * |
| H | 3 | * |
| J | 2 | * |
| L | 2 | * |
| R | 2 | * |
| S | 2 | * |
| U | 2 | * |

Third column of Table1 contains the sorted transactions of transaction database . Table (b) represents the fptree with header table with items sorted in the order of item names. In this section we introduce another way to construct a more efficient header tree by using another method. After the first scan we will get all frequent items with its frequency count. Sort the items in the order of its frequency and construct the header tree based on this list. The root of the header tree constructed by using this method contains the highest frequent item and all other items with highest frequency are stored around the root. The items with highest frequency will be searched more than less frequent items, so by using this method we can reduce the searching time again.

| Items | Frequency |
|-------|-----------|
| C | 6 |
| F | 5 |
| B | 4 |
| D | 4 |
| H | 3 |
| E | 3 |
| A | 2 |
| G | 2 |
| J | 2 |
| L | 2 |
| R | 2 |
| S | 2 |
| U | 2 |

III. RESULTS

The database contains different type of data transactions. Here we are giving review dataset as an input to the system for creating the frequent itemset to generate the association rule. The system uses hash table to store the item ID and the count of the frequency of the item so that when examining the candidate key, it only needs to check the hash table and this will make the process much efficient. After generating frequent itemsets for each length-n patterns, the system check transactions that do not contain any frequent n-itemsets and remove them from the itemset.

The Dataset can be uploaded to the system to perform the further operation

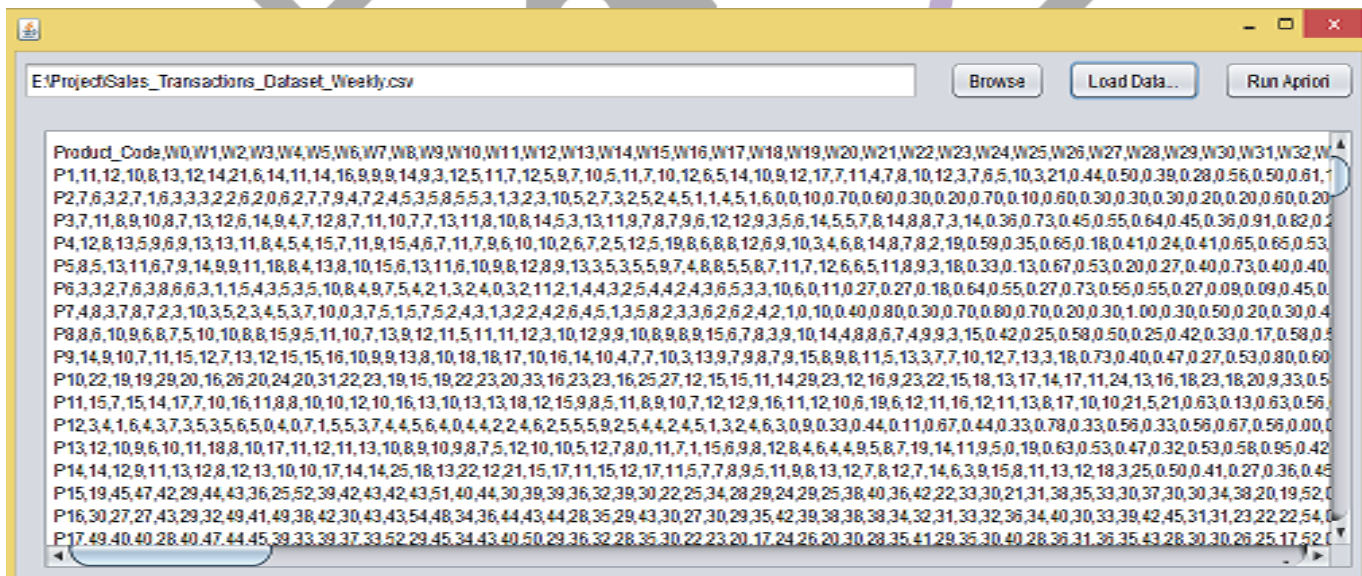


Fig. Dataset Upload

In the dataset upload in the application we will data from the system and load it to the application. The dataset should be comma separated CSV file. It can run Apriori algorithm, FP-Growth, Improved on the dataset.

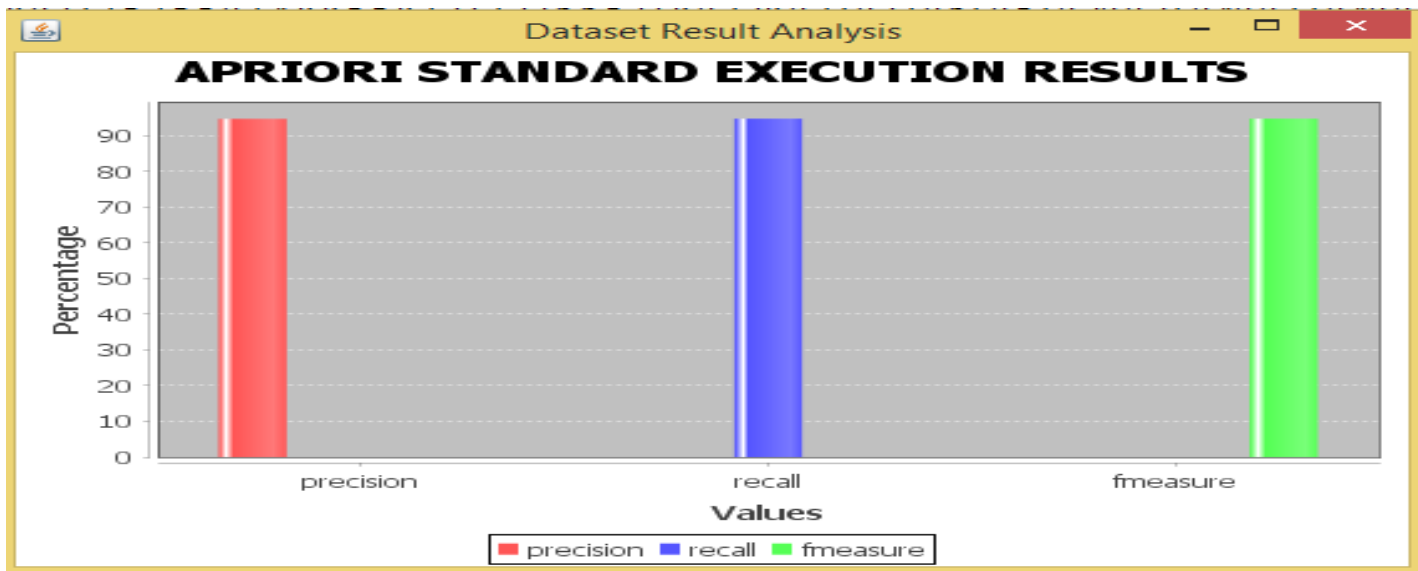


Fig: Standard Apriori Algorithm

The performs measure for standard algorithm the efficiency of the algorithm is based on the precision, recall , fmeasure. It uses the candidate key generation for calculating the frequent items from the dataset.

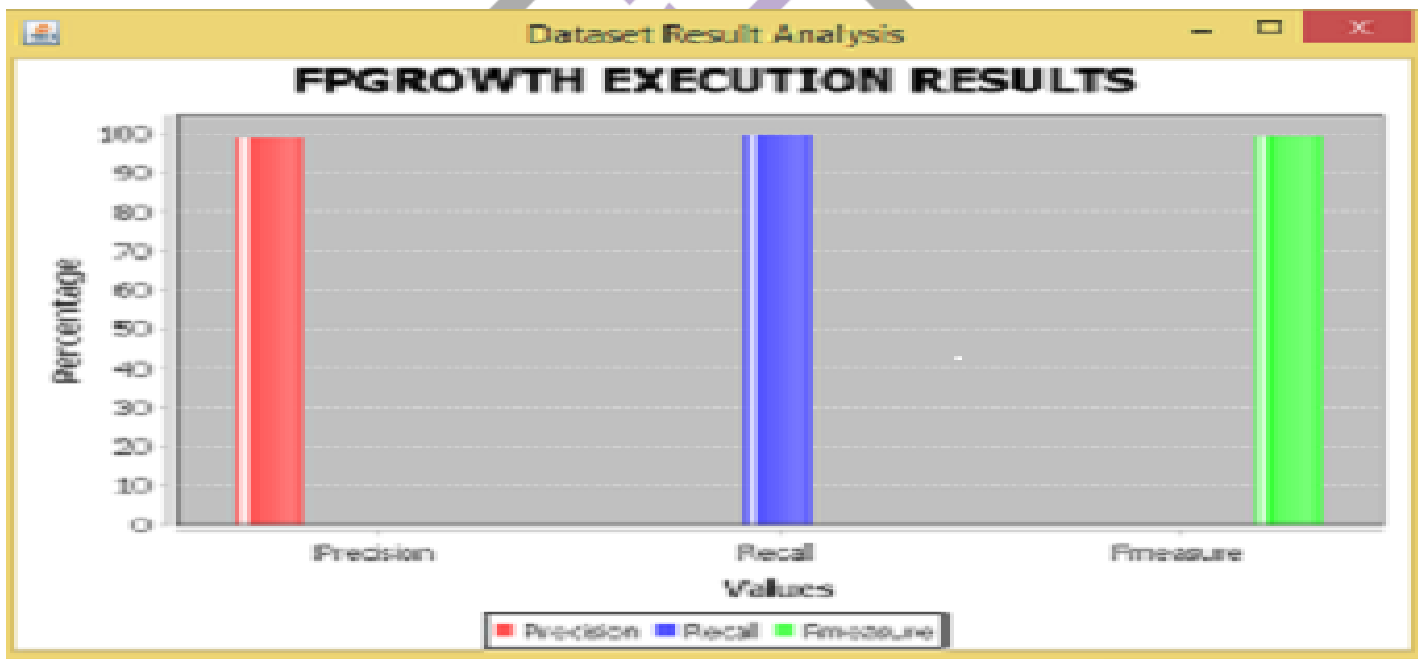


Fig:FP-Growth Algorithm

FP-Growth execution results are shown in fig. The efficiency of the FP-Growth algorithm is better than the standard Apriori algorithm. The precision of the algorithm is better than standard It will take less time.

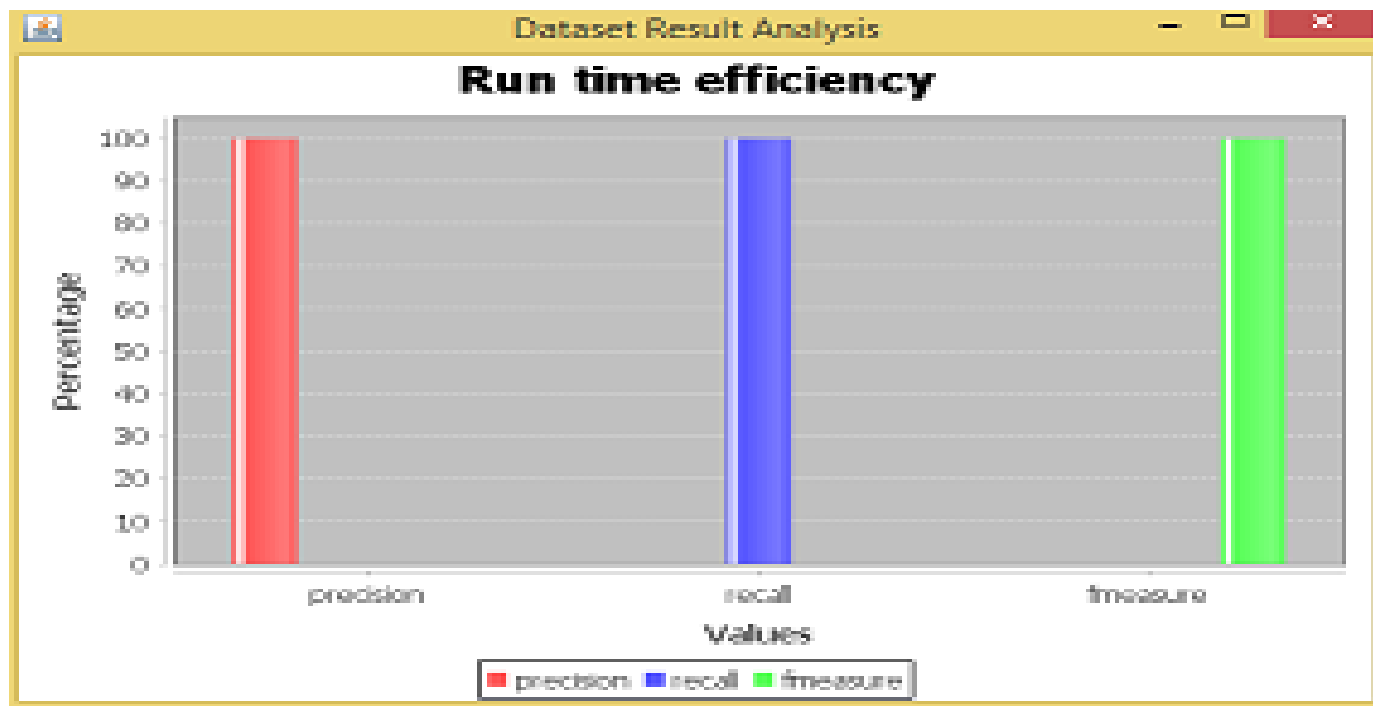


Fig.3 Improved Apriori Algorithm

As we upload dataset on the system it will execute system. As we can observe from above graph the run time efficiency of Proposed system. It is the better than existing system.

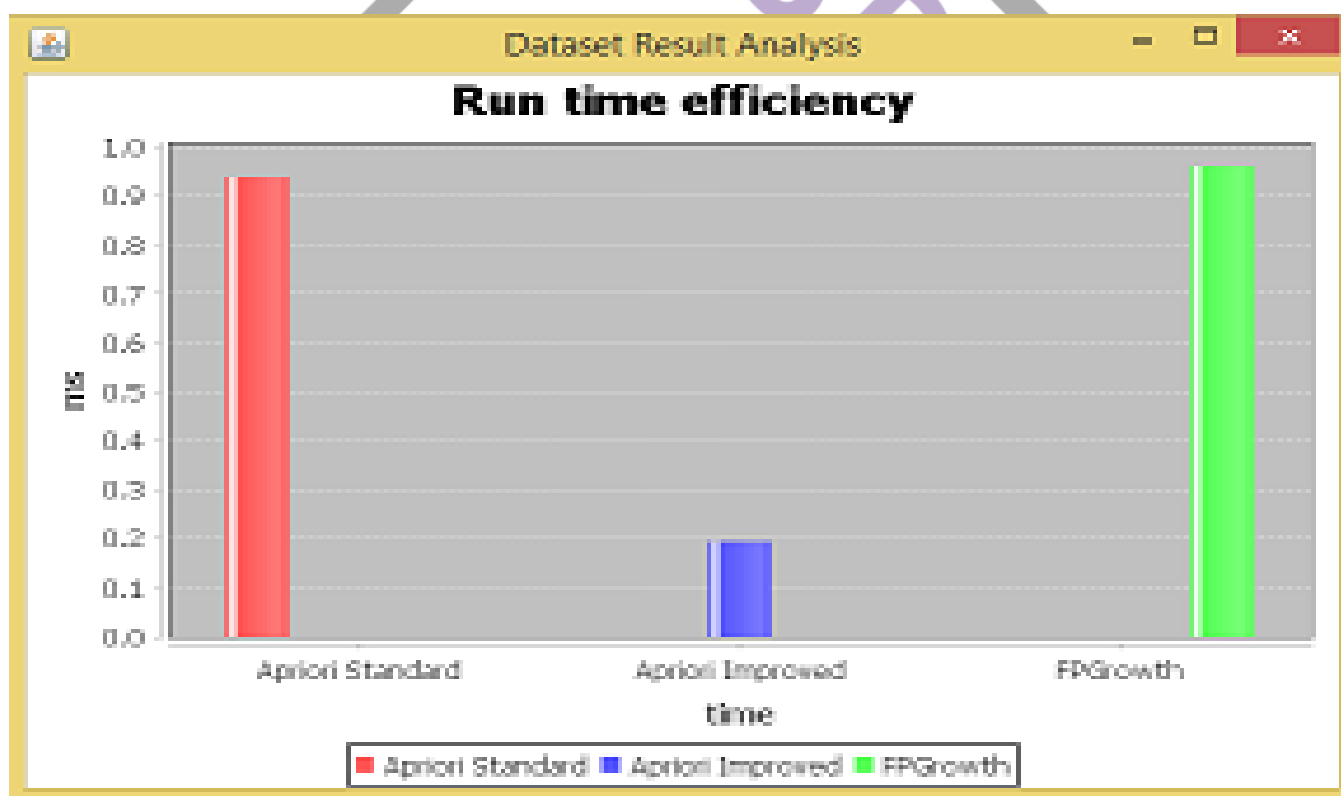


Fig.4 Comparison

The based approaches require more time for the existing system that are Standard Apriori Algorithm and FP-Growth Algorithm . The run time efficiency of proposed system is best.

IV. CONCLUSION

Frequent mining of database listings is important for the rules of mining associations. Frequent mining of database can be done by using the Apriori algorithm. Apriori uses a class method to create a model with one item, then two items, three items, and so on. In addition, it scans the database repeatedly to count the certifications of each form. In turn, FP Growth uses in-depth search

instead. Search for amplitude and use the pattern growth approach. To improve the performance of the Apriori algorithm, and reduce the time required for the sorting frequent itemset by using the binary search tree.

V. ACKNOWLEDGMENT

I have a tremendous pleasure in presenting the project “Mining Frequent Itemsets on Temporal Data By Using APriori Algorithm” under the guidance of Dr. D.V. Patil and PG coordinator Prof. A. S. Vaidya. I am really obligated and appreciative to Head of the Department Dr. D. V. Patil for their significant direction and consolation. I might likewise want to thank the Gokhale Education Society’s R. H. Sapat College Of Engineering, Management Studies Research, Nashik-5 for giving the required offices, Web get to and vital books. At last I must express my sincere heartfelt gratitude to all the Teaching Non-teaching Staff members of Computer Department of GESRHSCOE who helped me for their important time, support, remarks, thoughts.

REFERENCES

- [1] M. Ghorbani and M. Abessi “A New Methodology for Mining FrequentItemsets on Temporal Data” IEEE Trans. 2017
- [2] R. Agrawal, T. Imielinski, and A. Swami, “Mining association rules between sets of items in large databases,” ACM SIGMOD Rec., vol. 22, no. 2, pp. 207–216, 1993.S. Chaudhuri, V. Ganti, and R. Kaushik. A primitive operator for similarit joins in data cleaning. In Proc of ICDE’06, 2006.
- [3] R. Srikant and R. Agrawal, “Mining generalized association rules,” in Proc. 21st Int. Conf. Very Large Data Bases, 1995, pp. 407–419.
- [4] B. Liu, W. Hsu, and Y. Ma, “Integrating classification and association rule mining,” in Proc. 4th Int. Conf. Knowl. Discovery Data Mining, 1998, pp. 80–86.
- [5] M. Ale and G. H. Rossi, “An approach to discovering temporal association rules,” in Proc. ACM Symp. Appl. Comput.-vol. 1, 2000, pp. 294–300
- [6]] B. Liu, W. Hsu, and Y. Ma, “Integrating classification and association rule mining,” in Proc. 4th Int. Conf. Knowl. Discovery Data Mining, 1998, [4]S. Ramaswamy, S. Mahajan, and A. Silberschatz, “On the discovery of interesting patterns in association rules,” in Proc. 24th Int. Conf. Very Large Data Bases, 1998, pp. 368–379.
- [7] Y. Li, X. S. Wang, and S. Jajodia, “Discovering temporal patterns in multiple granularities,” in Temporal, Spatial, and Spatio-Temporal Data Mining. New York, NY, USA: Springer, 2001, pp. 5–19.
- [8] C.-H. Lee, C.-R. Lin, and M.-S. Chen, “On mining general temporal association rules in a publication database,” in Proc. IEEE Int. Conf. Data Mining, 2001, pp. 337–344.
- [9] C.-H. Lee, M.-S. Chen, and C.-R. Lin, “Progressive partition miner: An efficient algorithm for mining general temporal association rules,” IEEE Trans. Knowledge. Data Eng., vol. 15, no. 4, pp. 1004–1017, Jul./Aug. 2003.
- [10] X. Wu, C. Zhang, and S. Zhang, “Efficient mining of both positive and negative association rules,” ACM Trans. Inf. Syst., vol. 22, no. 3, pp. 381–405, 2004.
- [11] J.-W. Huang, B.-R. Dai, and M.-S. Chen, “Twain: Two-end association miner with precise frequent exhibition periods,” ACM Trans. Knowl. Discovery Data, vol. 1, no. 2, 2007, Art. no. 8.
- [12] J.-W. Huang, B.-R. Dai, and M.-S. Chen, “Twain: Two-end association miner with precise frequent exhibition periods,” ACM Trans. Knowl. Discovery Data, vol. 1, no. 2, 2007, Art. no. 8.
- [13] Y.-L. Chen and C.-H. Weng, “Mining fuzzy association rules from questionnaire data,” Knowl.-Based Syst., vol. 22, no. 1, pp. 46–56, 2009.
- [14] J. Adhikari and P. Rao, “Identifying calendar-based periodic patterns,” in Emerging Paradigms in Machine Learning. New York, NY, USA: Springer, 2013, pp. 329–357.
- [15] J. Adhikari and P. Rao, “Identifying calendar-based periodic patterns,” in Emerging Paradigms in Machine Learning. New York, NY, USA: Springer, 2013, pp. 329–357.
- [16] D. Nguyen, B. Vo, and B. Le, “CCAR: An efficient method for mining class association rules with itemset constraints,” Eng. Appl. Artif. Intell., vol. 37, pp. 115–124, 2015.