

# Improved Apriori with Sequential Pattern Mining Framework for Frequent Itemset Mining in Data mining

Pinky Sen<sup>1</sup>, Shubham Gangrade<sup>2</sup>, Manish Rai<sup>3</sup>

M.Tech. Scholar<sup>1</sup>, Guide<sup>2</sup>, Guide<sup>3</sup>

Department of Computer Science and Engineering, RKDF College of Engineering, Bhopal

**Abstract**—Data Mining is an analytic process designed for search of exact patterns and/or proper relationship between variables to scrutinize data and then to approve the exploration with the help of detected patterns and to new subsets of data. Frequent pattern mining is most easily explained by introducing *market basket analysis* (or affinity analysis), characteristic practice for which it is familiar. In this paper, an efficient hybrid algorithm was designed using a unifying process of the algorithms Improved Apriori and SPMF ECLAT. Results indicate that the proposed algorithm faster execution time and generate frequent rules with different-different support count.

**Keywords**—Data Mining, Hybrid Algorithm, Frequent Itemset, Improved Apriori, SPMF Eclat.

## I. INTRODUCTION

Data mining is one of these research activities to get familiar and comprehensible, before unidentified and convenient data from introductory data. This is finding the knowledge in database [1].

Data mining allows obtaining the important data from a large database. Data mining [2] is a procedure to discover stimulating information from a wide range of data put in storage in data mining, data maintenance, or additional data storing. In a outsized database of sales communications, relationship between objects is mining for laws the database has been a main part of study. From the analysis of a large Set of Super market [3], the real problem is defined by association rule mining was to search a correlation among sales of different products [4].

Association rule mining (ARM) [5] is the important part of data mining, which helps to predict the association among multiple data items. The big challenge of ARM is efficiently extract the knowledge from large size databases of various applications. As per concern of data holder, the main challenge of ARM is to share the accurate information with protection of sensitive information. To achieve this, Privacy preserving ARM plays very important role[6].

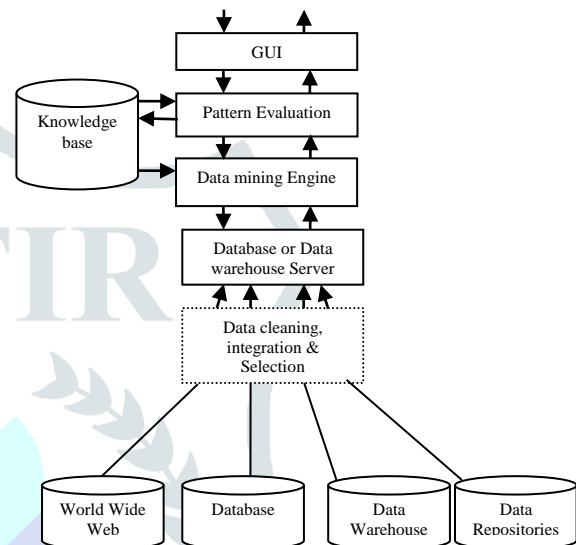


Fig. 1:Data Mining Architecture

In the rest paper, we have defined various literatures about related work in section II. In section III, described about proposed work that has been done. After this, we have discussed about experiment and results part in section IV. And in section V, we have concluded this paper.

## II. LITERATURE SURVEY

Yuxin Wang, TongkunXu, ShiqingXue, and YanmingShen [2018] proposed Dynamic queue & Deep Parallel Apriori (D2P-Apriori), a parallel frequent itemset mining algorithm on GPU to satisfy the high-performance requirement. The contributions of D2P-Apriori include as follows. The dynamic queue with bitmap is improved upon the vertical data structure to deal with the problem that the required memory for sparse data may exceed the size of GPU global memory. The Graph-join way is devised to adapt GPU architecture for candidate generation method. And the improved data structure also contributes to the significantly accelerated performance in support counting method on GPU. [8]

Said Jabboura, Fatima Ezzahra El Mazourib, LakhdarSaisa [2018] Using the controls, a negative approach has been suggested for strong negative association rules. The problem with this limitation shows us as a model. We have proven that we can benefit from very clear and strange SAT-based

approaches. Conviction interestingness is formulas as a nonlinear constraint and managed efficiently by choosing a particular branching heuristic. [9]

A. Naumoski, G. Mirceva, K. Trivodaliev & K. Mitreski[2018] Presented results in this paper, shows the influence of the three fuzzy metrics on the performance of the fuzzy pattern tree method used for learning diatom conservation representations. We experimentally evaluated, the membership functions, similarity metrics, model complexity and fuzzy operators on both train and test models. The results show what model parameters were best to set to obtain the highest accurate predictive model, after that two ecological models were learned, presented and interpreted. [10]

Yang Xi and Qi Yuan [2017] analyze the development status of intelligent tourist guide system and studies Apriori-based association rule algorithm with corresponding Development. This paper proceeds away the weight procedure of a subjective relationship. On this basis, we body a contained commendation model for the Tourists, and analyze instant data using data mining, to ensure more accurateness in the recommendation of the background, and to achieve more effective adapted spot recommend approaches.

Nurul Farise Sulkarnine, Ahmad Shah [2017] In this paper, the effectual hybrid procedure was considered with a united process of algorithm. and FP-Growth. The results show that despite the planned hybrid algorithm being more difficult, consume less memory resources & faster execution time [12].

HimaniBathla & Ms. Kavita Kathuria[2015] discussed various association rule algorithms and compared two algorithms: Apriori algorithm and Filter Associator. Through the Apriori algorithm, we have evaluated the frequent itemsets generation and number of cycle performed and also filter associator in the context of association analysis. According to the comparison of above two algorithms on weka tool, we conclude that Filter Associator is efficient algorithm than Apriori algorithm based on above two factors (Number of cycle performed, large itemsets) because the Apriori algorithm generates more number of cycle performed and generate extra large itemsets This reduces algorithm performance. [13]

Yong-le SUN and Ke-liang JIA [2009] find a new WSD method based on the mining association rules, also can coalfield the connotation instructions between the intelligence of the equivocal word and its context, and the sense of the equivocal word is find out by choosing the sense which results the most association rules concluded. The experimental results show that the precision obtained by our algorithm is higher than other algorithms. [14]

### III. PROPOSED METHODOLOGY

In the first part of the algorithm, the Improved Apriori property was used to discover all the maximal frequent item sets [15] which are repeating in the transactional database with a support value equal to or greater than the minimum support specified. There are still many itemsets which are frequent-1 but not included in the maximal frequent itemsets. So the database which contains frequent-1 elements are pruned but there are no maximal frequent itemsets which make the database smaller and easy to traverse. The pruned database

becomes the input in the second part of the algorithm which discovers all the frequent-1 itemsets and removes all the infrequent-1 itemsets from the transaction. Then, the SPMF Eclat algorithm was implemented from the pruned transactions.

#### A. Improved Apriori Algorithm

Apriori algorithm [16] is used for finding frequent itemsets in a dataset for Boolean association rule. Name of algorithm is Apriori is because it uses prior knowledge of frequent itemset properties. We use a real-time admittance or level visual search that uses k-frequent itemsets to discovery k + 1 items. An significant distinguishing of the atomic property we are observing for is that even diversities help to change the competence of a variety of substances. space.

Improved Apriori algorithm [17] removes the step to generate candidate itemsets which tend to improve the execution time of generating frequent itemsets. Improved Apriori depends on both forward and reverses scan of the given database. Uncertainty positive conditions are met, improved Apriori algorithm will decrease recurrence and scan time to discovery features in the package. This algorithm mines maximum frequent itemsets and their subset directly and makes a comparison with the items in the database. It prunes all the candidate itemsets according to the support count making sure that all the maximal frequent itemsets are mined.

#### B. Eclat Algorithm

Eclat [18] procedure is the first deep search algorithm. A perpendicular database builder is used in its place of a separate listing for all transactions; Each item is kept together in its cover (known as a tildel). If the substances are small in article numbers, the intracellular dishonorable method is used to calculate the support of an item less than apriori. It is suitable for small datasets and requires less time for frequent pattern generation than apriori.

#### Proposed Algorithm:

- Step:1 Input Retail dataset.
- Step:2 Scanning the transaction database one by one.
- Step:3 Create a 2-dimensional array; put the transactions at the count of repetition.
- Step:4 Arrange dataset in ascending order according to count of item in transaction.
- Step:5 Traverse the array to find maximal transactions. If there are no frequent itemset left go to step 7.
- Step:6 Take all the non-empty subset of a frequent itemset.
- Step:7 If there are frequent itemsets not included in the maximal itemset then find all the frequent-1 itemsets

and prune the database by removing the maximal frequent itemset.

Step:8 Find frequent-2 by the right neighbor method

Step:9 For Left Data We Use Eclat

Step:10 Save data as csv format

Step:11 Pre Process It using Unsupervised Numeric to Binary Attribute Filter.

Step:12 Choose SPMF Eclat from associate Tab.

1) Change data from horizontal layout to vertical layout. For both article incidences in changed communications.

2) It has multiplied a two-dimensional three-dimensional design. Adding TID\_sets of items in rows in a patterned perpendicular data format in the row pattern.

3) Using the new added items to generate new regular items using the existing frequency items in the pattern tree. Finally, all the corridors can be found by taking all the nodes of the pattern tree.

4) Advance information The entire contender is used to clip entirely in creating new regular items. Two items in creating an interest and calculating the support number. Take only which item (from pair) who have min-support

5) This process continues, for k+1 time, until no frequent items or no candidate itemsets can be found.

Step:13 Finally, we get all frequent itemset whose fulfill minimum support and threshold.

Step:14 Stop.

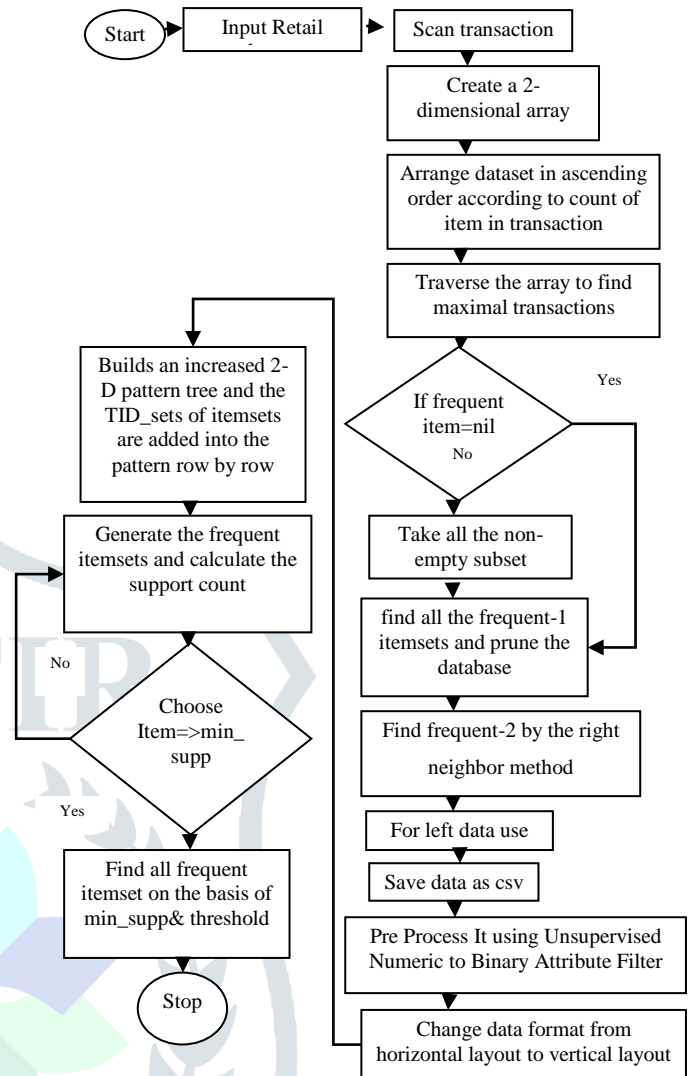


Fig. 3.1: Flowchart of Proposed Work

IV. EXPERIMENTAL RESULT & ANALYSIS

In the result analysis, the experiment of proposed work performed by using MATLAB [19] & WEKA [20]. Retail dataset 2006 is used for the investigational study of the frequent set mining. The resulting analysis shows the comparison of HYBRID Techniques with improved apriori& SPMF Eclat on selected datasets.

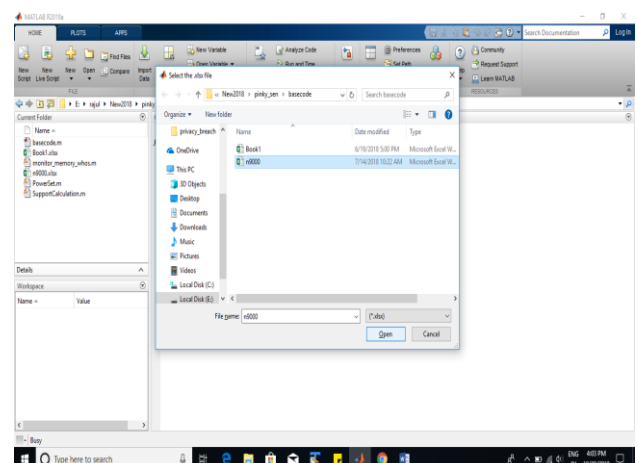


Fig 4.1: Select retail database

In fig 4.1 we have select the dataset (i. e. retail dataset). In this dataset considered total up 9000 transaction entries. This dataset is used for further process.

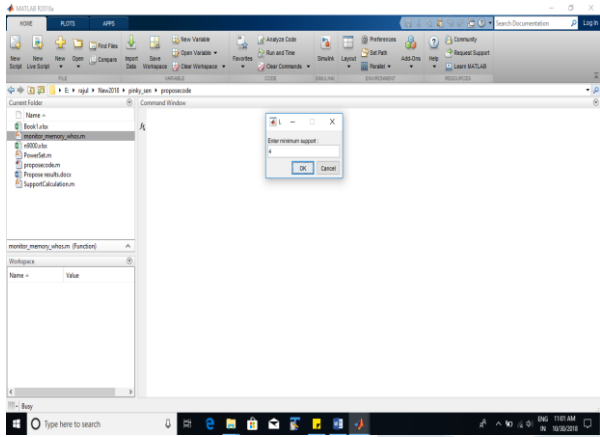


Fig 4.2: Input Support Count

In fig 4.2 shows entry dialog box for enter the value of minimum support, in which gives the minimum support value. We have given 4 as a minimum support value with Improved Apriori.

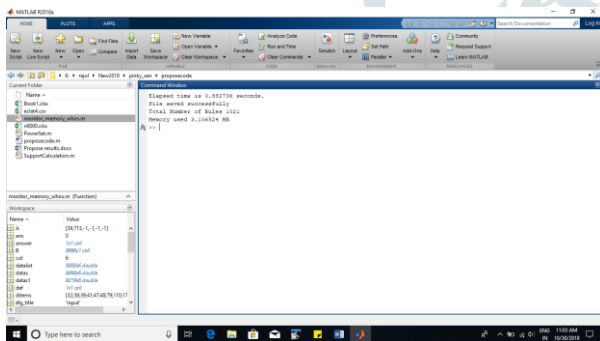


Fig 4.3: Result through improved apriori

For Left Data, We use Eclat algorithm by WEKA tool. For this select workbench from WEKA GUI displayed in fig. 4.4.

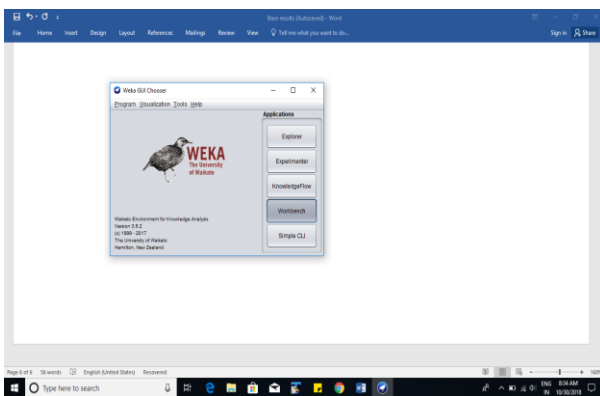


Fig 4.4:Select Workbench

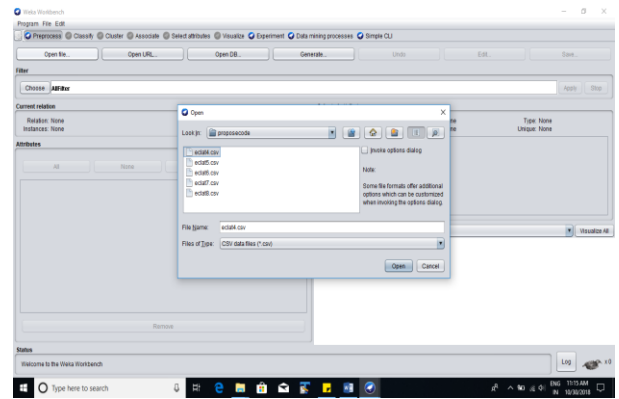


Fig 4.5:Open File

Pre processed the data using Unsupervised Numeric to Binary Attribute Filter in fig 4.6.

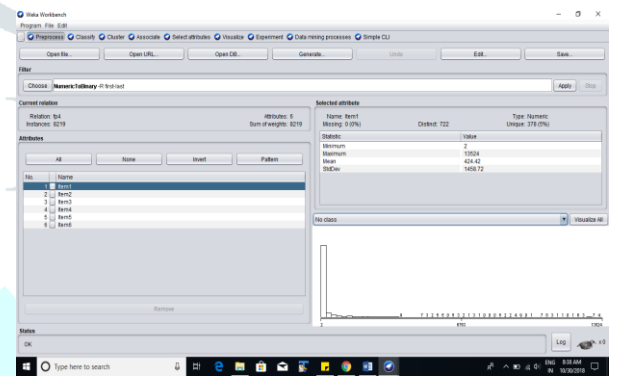


Fig 4.6:Pre Processing

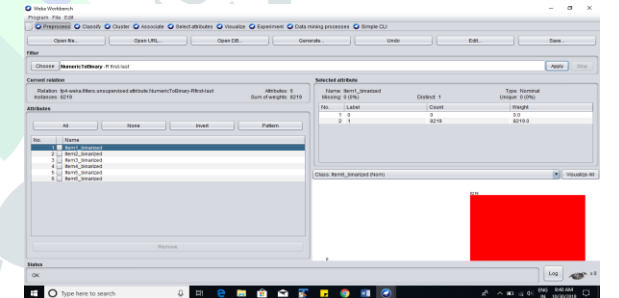


Fig 4.7:After Applying Filter

After performing preprocessing choose SPMF Eclat from associate Tab

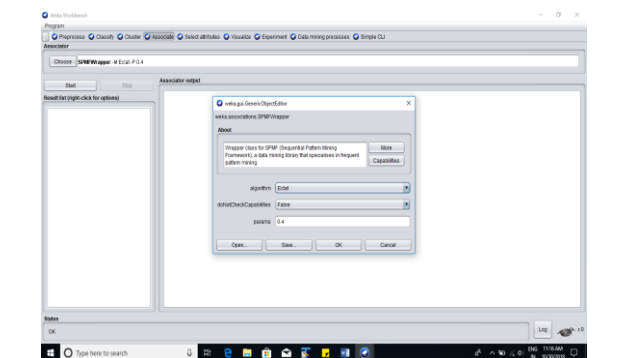


Fig 4.8:Choose SPMF Eclat from associate Tab



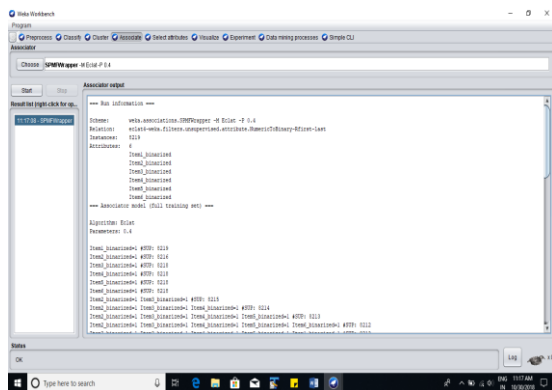


Fig 4.9: After Applying Eclat

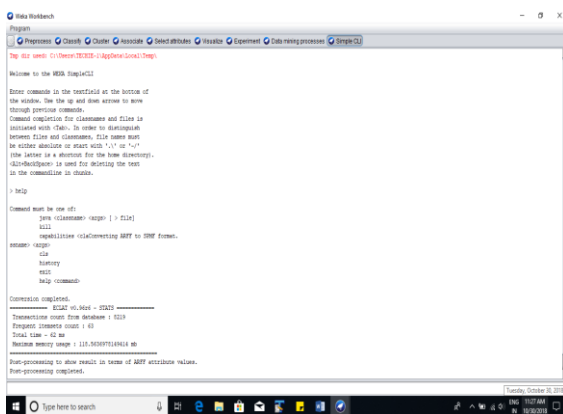


Fig 4.10: Post processing completed

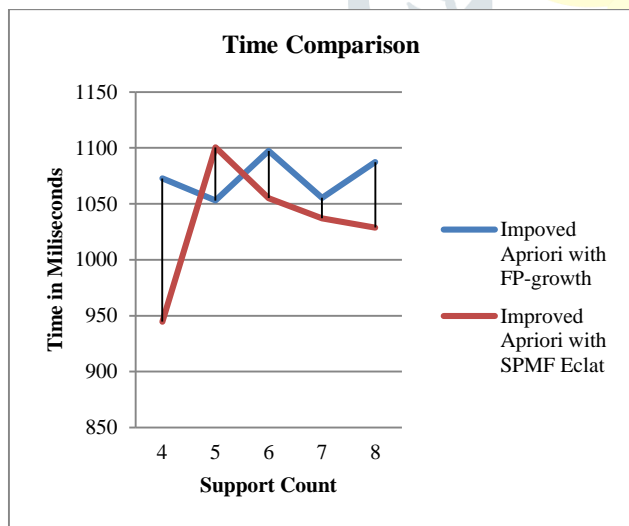


Fig 4.10: Elapsed time comparison

V. CONCLUSION

Now a day, data mining is used in almost all the places where a large amount of data is stored and processed. Suggestion Mining Explorations for Determined Items in the Data Set. Infrequently interacting in mines and stimulating relations and connections between the item sets in interconnected folders. In

this work, retail dataset has been used for mining the frequent item set. It is done by using MATLAB & WEKA tool. Here, we propose the discovery of frequent item sets from large transactional datasets by designing a new hybrid algorithm, which unifies the strengths of two existing algorithms. The Improved Apriori with SPMF Eclat algorithm shows that it is better in performance through execution time.

REFERENCES

- [1] [1]Lakshmi, B. N., & Raghunandhan, G. H., “ A conceptual overview of data mining”, National Conference on Innovations in Emerging Technology, February, 2011.pp.27-32.
- [2] [2] Charmi Mehta, “Basics of Data Mining: A Survey Paper”, International Journal of Trend in Research and Development, Volume 4(2), 2017, pp. 40-41.
- [3] [3]Mrs. R. R. Shelke, Dr. R. V. Dharaskar and Dr. V. M. Thakare, “Data Mining For Supermarket Sale Analysis Using Association Rule”, International Journal of Trend in Scientific Research and Development, Volume 1(4), 2017, pp. 179-183.
- [4] Demšar J, Zupan B (2013) Orange: Data Mining Fruitful and Fun - A Historical Perspective. Informatica 37:55–60.
- [5] Ms. J. Omana, Ms. S. Monika and Ms. B. Deepika, “Survey on Efficiency of Association Rule Mining Techniques”, International Journal of Computer Science and Mobile Computing, Vol.6 Issue.4, April- 2017, pg. 5-8.
- [6] Madhura Karanjikar And S. V. Kedar, “Secure Association Rule Mining Using Bi-Eclat Algorithm On Vertically Partitioned Databases”, International Conference On Intelligent Sustainable Systems (Iciss), pp. 176-181, 2017.
- [7] Ruhul, “Architecture of typical data mining system”, e-blog, March 12, 2011.
- [8] Yuxin Wang, Tongkun Xu, Shiqing Xue, and Yanming Shen, “D2P-Apriori: A deep parallel frequent itemset mining algorithm with dynamic queue”, ICACI, 2018 IEEE.
- [9] Said Jabboura, Fatima Ezzahra El Mazourib, Lakhdar Saisa, “Mining Negatives Association Rules Using Constraints”, Procedia Computer Science 127 (2018) 481–488.
- [10] A. Naumoski, G. Mirceva, K. Trivodaliev and K. Mitreski, “Learning Diatom Ecological Models with Fuzzy Order Data Mining Algorithm”, MIPRO, 2018.
- [11] Yang Xi, and Qi Yuan, “Intelligent Recommendation Scheme of Scenic Spots Based on Association Rule Mining Algorithm”, 2017, IEEE.
- [12] Nurul Fariza Zulkurnain and Ahmad Shah, “HYBRID: An Efficient Unifying Process to Mine Frequent Itemsets”, IEEE 3rd International Conference on Engineering Technologies and Social Sciences, 2017, pp. 1-5.
- [13] Himani Bathla, Ms. Kavita Kathuria(2015)”Association Rule Mining: Algorithms Used” IJCSMC, Vol. 4, Issue. 6, June 2015, pg. 271 – 277.
- [14] Yong-le SUN Association Rules and Ke-liang JIA, “Research of Word Sense Disambiguation Based on Mining”, 2009, IEEE.
- [15] Karam Gouda and Mohammed J. Zaki, “GenMax: An Efficient Algorithm for Mining Maximal Frequent Itemsets”, Data Mining and Knowledge Discovery, 11, 1–20, 2005.
- [16] <https://www.geeksforgeeks.org/apriori-algorithm/>
- [17] Wei, Y.Q., Yang, R.H. and Liu, P.Y., 2009, August. An improved Apriori algorithm for association rules of mining. In IT in Medicine & Education, 2009. ITIME'09. IEEE International Symposium on (Vol. 1, pp. 942-946). IEEE.
- [18] Siddhrajsinh Solanki and Neha Soni, “A Survey on Frequent Pattern Mining Methods Apriori, Eclat, FP growth”, International Journal of Computer Techniques, pp. 86-89.
- [19] <https://www.tutorialspoint.com/matlab/>
- [20] Dr. Sudhir B. Jagtap and Dr. Kodge B. G, “Census Data Mining and Data Analysis using WEKA”, International Conference in “Emerging Trends in Science, Technology and Management-2013, pp. 35-40.