# Evaluation of Frequent Itemset Mining Algorithms - priori and FP Growth

[1]Jismy Joseph,[2]Dr.G. Kesavaraj

*[1]PhD Research Scholar, [2]Professor and head*
*Department of Computer Science, Vivekanandha College of Arts and*
*Science for Women (Autonomous), Elayampalayam, Thiruchengode, Tamil Nadu, India*

[1]jismyjoseph2018@gmail.com
[2]Dr.Kesavaraj@vicas.org

*Abstract-* **Nowadays the Frequentitemset mining (FIM) is an essential task for retrieving frequently occurring patterns, correlation, events or association in a transactional database. Understanding of such frequent patterns helps to take substantial decisions in decisive situations. Multiple algorithms are proposed for finding such patterns, however the time and space complexity of these algorithms rapidly increases with number of items in a dataset. So it is necessary to analyze the efficiency of these algorithms by using different datasets. The aim of this paper is to evaluate theperformance of frequent itemset mining algorithms, Apriori   and Frequent Pattern (FP) growth by comparing their features. This study shows that the FP-growth algorithm is more efficient than the Apriori algorithm for generating rules and frequent pattern mining.**

*Keywords-* **Frequent itemset mining, Apriorialgorithm, FP-growth algorithm, Weka, Association Rule.**

## I.   INTRODUCTION

These days the size of databasesincreasesrapidly, this leads to invention of different tools to extract information automatically from the large database. Many researchers used data mining or knowledge discovery in database (KDD) to extract innovative and frequent pattern from large databases or transactional databases.

Frequent itemset mining has been applied in a great number of fields, including intrusion detection, Market basket analysis and credit card fraud prevention to discover unpredicted association among itemset in transactional and relational database.Frequent pattern Mining or association mining shows which items appear together in a transaction or relation.

Frequent patterns are patterns such as subsequences, itemsets or substructures that appears in a data set frequently. Finding such frequent pattern plays an essential role in mining association, correlations, and many other relationships among the data. It also helps in data classification, clustering and other data mining tasks as well. Thus, frequent pattern mining has becomes an important data mining task and a focused theme in data mining research [1].

## II.   LITERATURE REVIEW

Much research work has been done to compare different algorithms used for frequent itemset.  In the paper 'Survey on Frequent Item set Mining Algorithms' [2] Pramod S and O.P. Vyas conducted a survey on frquent itemset mining algorithms. They used Adult and Mushroom dataset for evaluating the performance of different algorithm. This study comprises features like different support values, size of transactions and different datasets. They concluded that SaM algorithm performed better in all data set.

In [3] M.S Mythili, A.R Mohamed Shanavas conducted a Performance Evaluation of Apriori and FP-Growth Algorithms by comparing the capabilities of these algorithms. The study show that FP-growth algorithm is more efficient than the Apriori algorithm.

In [4] Mr.Rahul Shukla, Dr. Anil kumar Solanki conducted a performance Evaluation for Frequent Pattern mining Algorithm. In the paper they compared the performance of Apriori and ECLAT Algorithm on medical data and they concluded that the Eclat approach is more efficient for mining frequent patterns in a large database.

In the paper 'Algoriihms for frequent itemset mining: a literature review' Chin-Hoong Chee, JafreezalJaafar, Izzatdin Abdul Aziz, MohdHilmi Hasan & William Yeoh reviewedthe strengths and weaknesses of the important and recent algorithms in Frequent Pattern Mining and they identified two major problem with frequent pattern mining. They mentioned that frequent hidden pattern mining needs more time and memory space [5].

In the paper[6]' Frequent Itemsets Mining for Big Data: A Comparative Analysis'DanieleApiletti,ElenaBaralis,TaniaCerquitelli,PaoloGarza,FabioPulvirenti and LucaVenturini conducted a theoretical and experimental comparative analyses of Hadoop- and Spark-based scalable algorithms to find frequent pattern from big data.

In [7] Ravi Ranjan and Aditi Sharmathey compared Hadoop, Spark, Flink by using Apriori and Fp-Growth on different dataset. They suggested that Flink is performing better in the field of big data.

## III. FREQUENT PATTERN MINING(FPM) AND ALGORITHMS

In data mining tasks, the frequent item sets plays an important role for finding frequent pattern or association from various kinds of databases like relational database, transactional database etc. FPM plays an essential role for clustering, classifying and identifying outliers a set of data. Apart from this, FPM has many applications like spatiotemporal data analysis, biological data analysis, and software bug detection [8].

FPM is used to predict the occurrence of a specific item based on the occurrence of other items in the transaction. The terminologies used in frequent pattern mining are support and confidence. Support specifies how frequently an itemset appears in the dataset and the confidence describes how often the rule has been found to be true. The support and confidence are defined as:

$$Support\ (A \Rightarrow B) = P\ (A \cup B)$$

$$Confidence\ (A \Rightarrow B) = \frac{Support\ (A \cup B)}{Support\ (A)}$$

The rule is considered as strong if it satisfies a minimum support threshold and a minimum confidence threshold. The methods used in frequent item set mining are

- Finding Frequent itemsets using candidate generation
- Mining Frequent itemsets without candidate generation
- Mining Frequent itemsets using Vertical Data format.

### A. APRIORI Algorithm

R.Ararwal and R.Srikant proposed this algorithm in 1994 for mining frequent itemsets for Boolean association rule [1]. It uses candidate generation for finding frequent itemsets. Aprioriis an iterative approach, where k-itemsets are used to find (k=1)-itemsets. This algorithm uses Apriori property to reduce the search space. It is a two-step process called join and prune.

*Join step* - To find $L_k$, a set of k-itemset is generated by joing $L_{k-1}$ with itself.

*Prune step* - Scans the count of each item. If it less than minimum support then it does not considered as frequent.

*Apriori Algorithm Pseudocode*

```
Procedure Apriori(D, minSupport)   // D – database,
minsupport- Minimum support
{
L1 = {frequent items};
For (k= 2; Lk-1! =∅; k++)
{
        Ck= candidates generated from Lk-1
        For each transaction t in database do
        {
                Increment all candidates
                Lk = candidates in Ck with
                minSupport
        }
}
ReturnCkLk;
}
```

### B. FP-Growth Algorithm

The main disadvantages of Apriori algorithm are it generates huge number of candidate sets and this algorithm repeatedly scans the database and checks a large number of candidates by pattern matching. Hence it is very costly. To overcome these

disadvantages, the next method for generating frequent itemset without using candidate generation is FP-Growth. FP-Growth uses a divide and conquer approach. It first compresses the database into FP-Tree and then divides these compressed database into a set of conditional database.

*FP-Growth Algorithm*

Step1: Scan the database to get set of frequent itemsetand their support count.

Step2: Sort the frequent itemset in descending order using support count.

Step3: Construct FP-tree. Initially it creates the root of the tree and labelled as 'Null'.

Step4: Construct the FP-conditional tree for each item (or itemset)

Step5: Determine the frequent patterns.

## IV.     DATASETS AND RESULTS

We have used supermarket data set and vote datasets from 'storm.cis.fordham.edu' for comparing two frequent pattern mining algorithms. The first data set is 'Super Market Data Set (SMDS)' [9], which contains 4627 instances and 217attributes. The second one is 'Vote Data Set (VDS)' [9], which consists of 435 instances and 17 attributes. In this comparative study, twofrequent itemset mining algorithms are used. They are Apriori and FP growth.

Initially, the datasets are preprocessed and afterthat the algorithms are applied. The following figures shows the results obtained from weka after preprocessing. Fig.1 shows the result obtained after preprocessing the dataset 'SMDS.
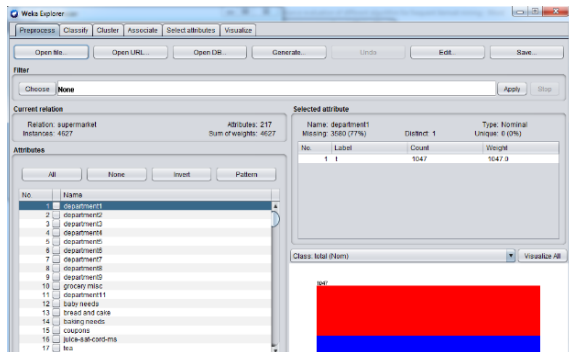

FIG.1. PREPROCESSING RESULT – SMDS

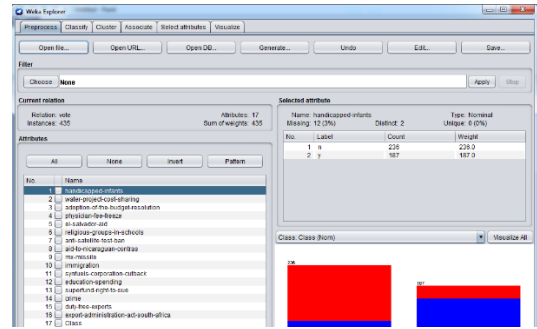Fig.2 shows the result obtained after preprocessing the dataset 'VDS'.


FIG.II. PREPROCESSING RESULT – VDS

The following figures shows the output obtained from Apriori and FP growth algorithms.Fig3 shows the rules generated by Apriori algorithm in SMDC dataset. Ten rules are generated based on support and confidence.
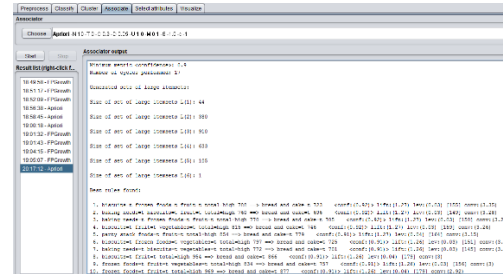

FIG.III.     THE     RULES     GENERATED     BY     APRIORI ALGORITHM IN SMDC DATASET

Fig.4 shows the rules generated by Apriori algorithm in VDS. Ten rules are generated based on support and confidence.
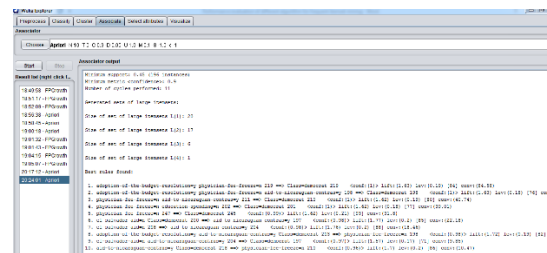

FIG.IV.     THE     RULES     GENERATED     BY     APRIORI ALGORITHM IN VDS

The rule generated by FP-Growth using vote data set is shown in Fig 5.  FP Growth found 41 rules in vote dataset and only the top 10 rules are displayed in the following figure.
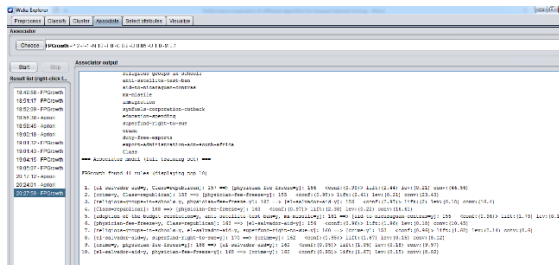
FIG.V. THE RULE GENERATED BY FP-GROWTH USING VOTE DATA SET

Fig 6 shows the rules generated by FP Growth in SMDC data set. FP Growth found 16 rules in vote data set and only the top 10 rules are displayed in the following figure.
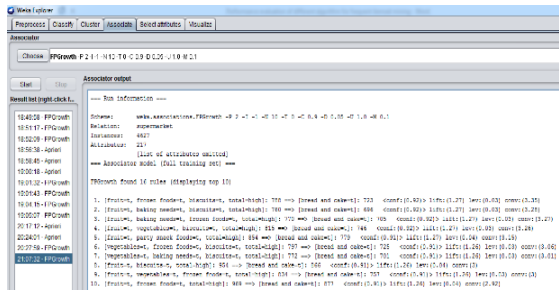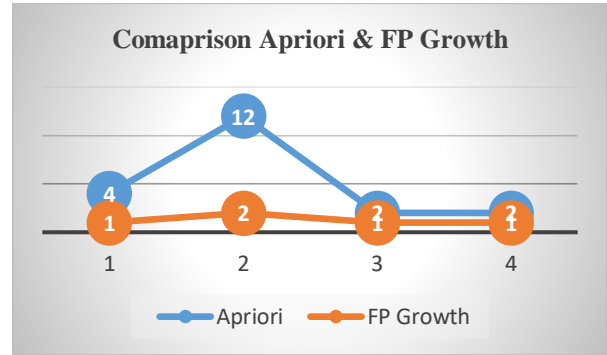


FIG. VI. THE RULES GENERATED BY FP GROWTH IN SMDC DATA SET

## V. COMPARISON OF RESULTS

This section deals the comparison part of Apriori and FP Growth algorithms. The table1 shows the time needed to generates rules by using Apriori and FP Growth.Compared to Apriori the FP Growth algorithm takes less time to generate the rules in both datasets. In the supermarket dataset,Apriori algorithm took 4 seconds to produce the rules when minimum support is 0.15. But FP growth algorithm generated the rules within in a second. In the case of vote dataset FP Growth performed better than the Apriori algorithm.

| Instances | Attribute | Minimum support | Time Taken (in Sec) | |
|---|---|---|---|---|
| | | | Apriori | FP Growth |
| SupperMarket (4627) | 217 | 0.15 | 4 | 1 |
| | | 0.2 | 12 | 2 |
| Vote(435) | 17 | 0.15 | 2 | 1 |
| | | 0.2 | 2 | 1 |

TABLE 1- COMPARISON RESULTS



Comaprison Apriori & FP Growth

## VI. CONCLUSION

In many data mining applications, association rule plays an important role for finding frequent pattern. In this study we observed that FP Growth algorithm is better than the Apriori algorithm. In both datasets the FP growth taken less time to generate the rule. FP-growth is more acceptable for larger databases.

*References*

[1] Jiawei Han and Micheline Kamber. Data Mining: Concepts and Techniques, Morgan Kaufmann Publishers, Second Edition, 2003.

[2] Pramod, S., & Vyas, O. P. (2010). Survey on frequent item set mining algorithms. *International journal of computer applications*, *1*(15), 86-91.

[3] Mythili, M. S., &Shanavas, A. M. (2013). Performance evaluation of apriori and fp-growth algorithms. *International Journal of Computer Applications*, *79*(10).

[4] Mr.Rahul Shukla, &Dr.AnilkumarSolanki(2015). Performance Evaluation for Frequent Pattern mining Algorithm. International Journal of Engineering Research and General Science Volume 3, Issue 5.

[5] Chee, C. H., Jaafar, J., Aziz, I. A., Hasan, M. H., &Yeoh, W. (2019). Algorithms for frequent itemset mining: a literature review. *Artificial Intelligence Review*, *52*(4), 2603-2621.

[6] Apiletti, D., Baralis, E., Cerquitelli, T., Garza, P., Pulvirenti, F., &Venturini, L. (2017). Frequent itemsets mining for big data: a comparative analysis. *Big Data Research*, *9*, 67-83.

[7] Ranjan, R., & Sharma, A. (2019). Evaluation of frequent itemset mining platforms using apriori and fp-growth algorithm. *International Journal of Information Systems & Management Science*, *2*(2).

[8] Aggarwal CC (2014) An introduction to Frequent Pattern Mining. In: Aggarwal CC, Han J (eds) Frequent Pattern Mining. Springer, Basel, pp 1–14.

[9] Weiss, G. M. (2020, September). *data-mining/datasets.html*. Retrieved from storm.cis.fordham.edu: https://storm.cis.fordham.edu/~gweiss/data-mining/datasets.html