# Frequent Item set Mining Using Global Profit Weight Algorithm

ASHA RAJKUMAR M.phil, Reasearch Scholar
Computer Science Department
P.S.G.R.Krishnammal College for Women
Coimbatore, India

G.SOPHIA REENA M.C.A., Mphil.,
HOD, BCA Department
P.S.G.R.Krishnammal College for Women
Coimbatore, India

*Abstract*— The objective of the study focused on weighted based frequent item set mining. The base paper has proposed multi criteria based frequent item set for weight calculation. Contribution towards this project is to implement the global profit weight measure and test the performance over utility based mining. For this project the data consist of 90 products from automobile shop including unit price, quantity sold and profit margin for transaction set (one month data). Algorithm has been implemented in Visual Basic for visualizing step by step process calculations**.** Supervised machine learning techniques namely Naïve Bayes Decision tree classifier, VFI and IB1 Classifier are used for learning the model. The results of the models are compared and observed that Naïve Bayes performs well. WEKA tool is used to classify the data set and accuracy is calculated.

Keywords- Global Profit Weight Algorithm; Classification Algorithm; WEKA Tool

## I. INTRODUCTION

Discovery of efficient association rules has been found useful in many applications. However, without fully considering the importance and significance of items and transactions, it is noted that some discovered rules may be expired from users' interest. Utility measures play an important role in data mining, regardless of the kind of patterns being mined. These measures are intended for selecting and ranking patterns according to their potential interest to the user. This study focused to implement the Global Profit Weight (GPW) for the frequent item set. Generally the profit can be measured in a traditional way. In this study they proposed multi criteria based profit calculation.

This work employs global profit weight algorithm is implemented using visual basic to find the profit of the item set in the transaction. Classification algorithm is used to evaluate the profit measure such as high (H), medium (M) and low (L). Widely used supervised machine learning techniques namely Naïve Bayes Decision tree classifier, VFI and IB1 Classifier are used for learning the model. The results of the models are compared and observed that Naïve Bayes performs well. WEKA tool is used to classify the data set and accuracy is calculated.

## II. FREQUENT ITEM SET MINING

Frequent item sets may only contribute a small portion of the overall profit, whereas non-frequent item sets may contribute a large portion of the profit. In reality, a retail business may be interested in identifying its most valuable customers (customers who contribute a major fraction of the profits to the company). Hence, frequency is not sufficient to answer questions, such as whether an item set is highly profitable, or whether an item set has a strong impact. Profit mining is thus useful in a wide range of practical applications and was recently studied in [15].

### A. Analytical Hierachical Process

This paper presents Analytic Hierarchy Process and it uses pair wise comparisons and computes the weighting factors. The utility measure [6], was proposed to overcome the shortcomings of support. It reflects the impact of the quality sold on the cost or profit of an item set. Lu et al [12], Proposed a scheme for weighting each item using a constant value without regard to the signification of transactions [12]. In this scheme, the utilities are attached to the items rather than the transactions. Wang et al [14] suggested that it remains unclear to what extent patterns can be used to maximize the business profit for an enterprise.

The project generalizes previous work on profit measure. The profit measure of the items in the transaction is defined by using the characteristic of the item. Considering the profit of an item, there are a number of important factors to consider as well. This paper defines five variables for items to compute the Global Profit Weight. (i.e. Damage, Offer, Quality, Margin and Frequency). The preference of the Quality factor is calculated according to the level of Quality measures.

### III. ASSOCIATION RULE

Association Rule is an important type of knowledge representation revealing implicit relationships among the items present in large number of transactions. Given $I=\{i_1,i_2,i_3,\ldots,i_n\}$ as the items' space, which is a set of items, a transaction may be defined as a subset of I, and a dataset may therefore be defined as a set D of transactions. X and Y are non-empty subsets of I. The support of an item set X in a dataset D, denoted as support D(X), is defined as count D(X)/|D|, where count D(X) is the number of transactions in D containing X. An item set

is said to be frequent (large) if its support is larger than a user-specified value (also called minimum support (min_sup)). An association is an implication of the form [X → Y, sup, conf], where $X \subset I$, $Y \subset I$, and $X \cap Y = \emptyset$. The support of $X \cup Y$ (sup) in the transactions is larger than min_sup, furthermore when X appears in a transaction; Y is likely to appear in the same transaction with a probability conf. Given a threshold of minimum support and confidence, methods of discovering association rules [4, 5, 6, 9] have become active research topics since the publication of Agarwal, Imielinski and Swami and Agarwal and Srikant papers [2, 3].

However, the traditional Association Rule Mining (ARM) model assumes that items have the same significance without taking account of their weight/attributes within a transaction or within the whole item space. But this is not always the case. For example, [wine, salmon, 1%, 80%] may be more important than [bread, milk, 3%, 80%] even though the former holds a lower support. This is because those items in the first rule usually come with more profit per unit sale, but the standard ARM simply ignores this difference.

B    Problem Classification

The mining of association rules for the unweighted case has been done for several years. However, for the above reasons, association rules have been developed for weighted items. To begin with, the association rule must be defined first. Similar to [1], [2] consider a database with a set of transactions D, a set of attributes or items T, and each transaction is assigned a transaction identifier <TID>.

In this study, the major problem is to mine the association rules with weighted items, based on the different types of the association rules, which are binary association rules and quantitative association rules. New algorithms are required to solve such problems since the available algorithms cannot be solved.

## IV. ASSOCIATION RULE MINING

The search space of all association rules contains exactly $3^{|I|}$ different rules. However, given all frequent item sets, this search space immediately shrinks vastly. Indeed, for every frequent item set I, there exists at most $2^{|I|}$ rules of the form X⇒Y, such that $X \cap Y = I$. Again, in order to efficiently traverse this search space, sets of candidate association rules are iteratively generated and evaluated, until all frequent and confident association rules are found. The underlying technique to do this is based on a similar monotonicity property as was used for mining all frequent item sets.

Proposition 1. (Confidence monotonicity) Let X, Y,Z $\subseteq$I be three item sets, such that $X \cap Y = \{\}$. Then,
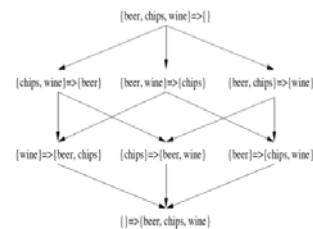


Figure 1.    Example of Association Rule.

Association rules (ARs) have been widely used to determine customer buying patterns from market basket data. The task of mining association rules is mainly to discover association rules (with strong support and high confidence) in large databases. Classical Association Rule Mining (CARM) deals with the relationships among the items present in transactional databases. The typical approach is to first generate all large (frequent) item sets (attribute sets) from which the set of ARs is derived. A large item set is defined as one that occurs more frequently in the given data set than a user supplied support threshold. To limit the number of ARs generated a confidence threshold is used. The number of ARs generated can therefore be influence by careful selection of the support and confidence thresholds, however great care must be taken to ensure that item sets with low support, but from which high confidence rules may be generated, are not omitted.

$$confidence(X \setminus Z \Rightarrow Y \cup Z) \leq confidence(X \Rightarrow Y).$$

*Proof.* Since $X \cup Y \subseteq X \cup Y \cup Z$, and $X \setminus Z \subseteq X$, we have

$$\frac{support(X \cup Y \cup Z)}{support(X \setminus Z)} \leq \frac{support(X \cup Y)}{support(X)}.$$

C    Weighted Association Rule Mining

One possible problem with the definition is that when the number of items in an item set is large, then the total weight may be large, even if each item has a small weight. In this section, we focus on the mining of weighted association rules for which the weight of an item set is normalized by the size of the item set. The choice of using unorganized or normalized weight will depend on the individual need of each application.

Although the semantics of the rules will be different, previous algorithm MINWAL(O) can be applied for this case, with a modification of the definitions of large weighted item sets and k-support bound However, MINWAL(W) new algorithm is applied.

D    Frequent Pattern  Mining Problem

Definition 10**.** Let D be a transaction database over a set of items I, and σ a minimal support threshold. The collection of frequent item sets in D with respect to σ is denoted by

$$\mathcal{F}(\mathcal{D},\sigma) := \{X \subseteq \mathcal{I} \mid support(X,\mathcal{D}) \geq \sigma\}$$

Example 1. Consider the database shown in Table 1 over the set of items I = {beer, chips, pizza, and wine}.

| Tid | X |
|-----|---|
| 100 | {beer, chips, wine} |
| 200 | {beer, chips} |
| 300 | {pizza, wine} |
| 500 | {chips, pizza} |

Table 1: An Example Transaction Database d.

Table 1 shows all frequent item sets in D with respect to a minimal support threshold of 1. Table 3 shows all frequent and confident association rules with a support threshold of 1 and a confidence threshold of 50%.

## V. MULTI CRITERIA DECISION MAKING

E    Analytical Hierachical Process

This process uses pair wise comparisons and then computes the weighting factors through evaluation of a set of criteria elements. The decision maker starts by laying out the overall hierarchy of the decision. This hierarchy reveals the factors to be considered as well as the various alternatives in the decision. Then a number of pair wise comparisons are done, which result in the determination of factor weights and factor evaluations [18]. The process has been used to assist numerous corporate and government decision makers.

In this process the problems are decomposed into a hierarchy of criteria and alternatives. An important part of the process is accomplished by the three steps. The steps are stating the Objective, Defining criteria, and Pick the alternatives. This information is then arranged in a hierarchical tree and synthesized to determine relative rankings of alternatives. Both qualitative and quantitative criteria can be compared using informed judgments to derive weights and priorities

For Example, We have three mobiles Nokia, Motorola and Sony. While Comparing Nokia with Motorola, Nokia slightly favored then Motorola, thus we put 1/3 in the row 1 column 2 of the matrix. Comparing Nokia and Sony, Nokia strongly more preferred (likes) than Sony, thus we put actual judgment 5 on the first row, last column of the matrix. Comparing Motorola and Sony, Motorola is very strongly preferred. Thus we put his actual judgment 7 on the second row, last column of the matrix.  Because we have three comparisons, thus we have 3 by 3 matrix. The diagonal elements of the matrix are always 1 and we only need to fill up the upper triangular matrix. To fill the lower triangular matrix, we use the reciprocal values of the upper diagonal. If aij is the element of row ith column jth of the matrix, then the lower diagonal is filled using this formula

$$a_{ji} = 1/a_{ij}$$

$$A = \begin{array}{c} \text{Nokia} \\ \text{Motorola} \\ \text{Sony} \end{array} \begin{bmatrix} 1 & 1/3 & 5 \\ 3 & 1 & 7 \\ 1/5 & 1/7 & 1 \end{bmatrix}$$

Fig 2: Pair wise Matrix

F    Global Profit Weight Algorithm

In this section we propose new GPW Algorithm, which is used to derive the list of Global Profit Weight for the set of items. The inputs are product (P), Quantity (Q), Unit Price (UP). The output derived is called as Global Profit Weight (GPW).

Algorithm          :          GPW
Input              :          Product (P) , Quantity (Q), Unit Price (UP)
Output             :          The set of Global Profit Weight (GPW)

Procedure Global Profit Weight (P, Q, UP)
        Begin
        GTP ← Ø;
        For each Pi ≤ n do
                        TPi ← Qi * UPi;
        GTP ← GTP + TPi;
        End
For each Pi ≤ n do
        TPPi ← TPi / GTPi;
End
Cwr ← ((∏Cr)1/Cn) / (Σ(∏Cr)1/Cn)
Pwk ← ((∏CPk)1/CPn)/ (Σ(∏CPk)1/CPn)
For each i ≤ n
   PWi ← Ø;
   For each x ≤ r
        For each y ≤ k
                PWi ← PWi + (Pwxy * Cwyx) ;
        End
End
For each Pi ≤ n do
GPWi ← PWi + TPPi;
End
End

G    GPW IMPLEMENTATION

The implementation is achieved with a real time transaction set collected from automobile shop.

Fig 3: GPW Model



Fig 4: Transaction Set

The above screen shows the transaction set collected for this study. The original is kept except the product ID. When a GPW model option is selected it shows the step 1 screen. Open data button open the data from file and display it in a grid.

The following Figure 5 shows the profit weight calculation from the transaction set. As discussed in the base paper, the calculation is made on total profit and profit weight.
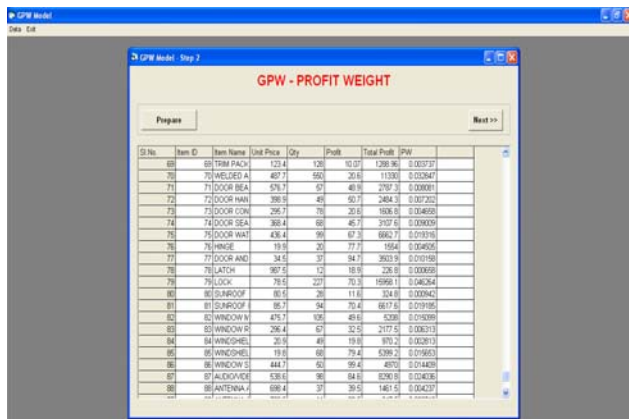


Fig 5: Profit Weight (PW)

As name described, the unit price shows the item price, qty column denotes number of units sold in one month, profit is marginal profit per item collected from the organization. Total profit is calculated through multiplying quantity and profit margin. The profit weight is calculated from total profit for the item divided by total profit earned from all the items.

The following Figure 6 describes the criteria selection. The criteria are built according to our objective of the study.
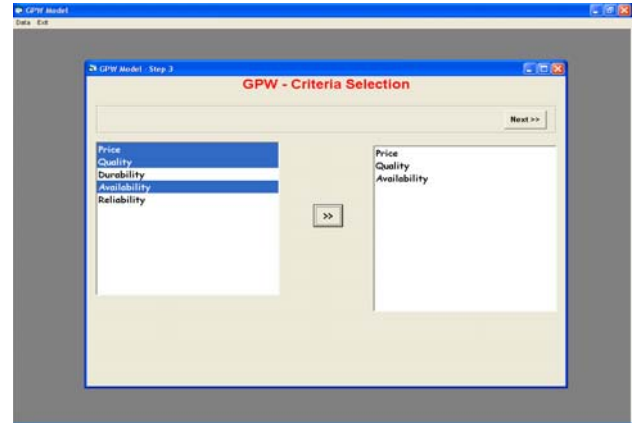


Fig 6: Criteria Selection

The main objective of our study is to predict the global profit weight with its importance measured through criteria. Here, the weight can be measured up to five criterions such as price, quality, durability, availability and reliability. But in our experiment we would like to measure the products with price, quality and availability band

The following Figure 7 shows the AHP calculation for criteria weight.



Fig 7: Criteria Weight

As shown in Figure 7, we have chosen only three criteria. In this screen we would like to calculate the criteria weight through AHP pair wise matrix. We have applied the same rule proposed in AHP nine point scaling technique. The p value is product weight and the z value is a priority ratio.

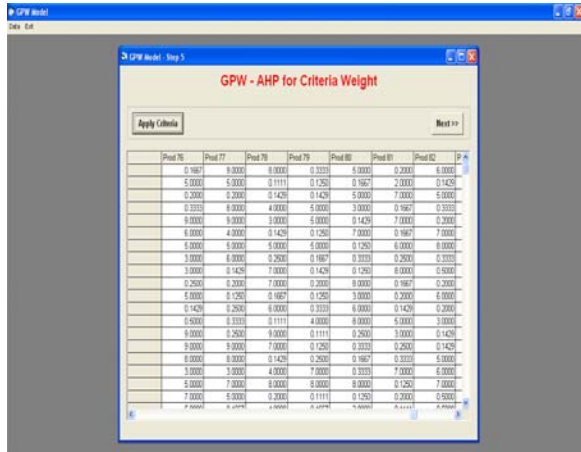The following Figure 8 shows the criteria based product weight.

Fig 8: Criteria Based Product Weight

In this screen it shows three matrixes with twenty five products each and also calculated the p value and z value. This screen describes that it is calculated for quality wise product priority, price wise product priority and availability wise product priority.

The following Figure 9 shows the global profit weight measure.



Fig 9: Global Profit Weight

It is inferred from the above screen that criteria wise product priority is shown in first three columns. The fourth column denotes the general criteria weight calculated in Figure 8. The Profit Criteria Weight (PCW) is computed as matrix multiplication of criteria wise product priorities and criteria weight. The global profit weight measure is an abstraction of PCW and Profit Weight (Figure 5).

The global profit weight is built up with two distinct values; marginal profit earned through sales and sentimental attribute values for all products. It helps to decide the best item in both senses. The product priority can be checked with the strategy as higher the value (GPW) is the best product.

H    Performance Evaluation

In the previous section describes about the model developed for global profit weight. The objective of GPW is an alternative for the traditional weighted association rule mining. An efficient algorithm is required because a significant amount of processing is undertaken to the application of weighted association rule mining. Experiment is performed with the real time data set collected from the automobile company. The GPW algorithm was implemented and presented in step wise manner.
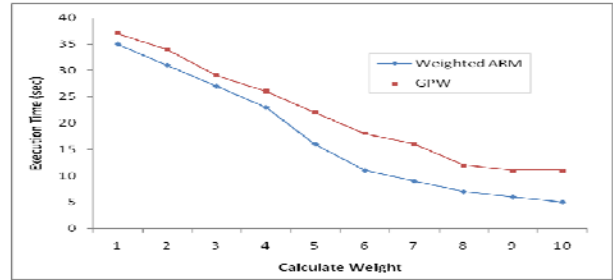


Fig 10: Execution time to generate weight for products

The above figure exhibits that the traditional weighted association rule mining algorithm has taken lesser time than global profit weight measure. The experiment has taken for 10 items. It is also noted that there is some similarity found between WARM and GPW on its execution time. The proposed algorithm is required high computation power to predict the weight measure.

To measure the product weight with its quality of output for 10 products. The quality is measured in five-point scaling like very poor quality, poor quality, fair quality, good quality and very good quality. The following figure 11 exhibits the scaling measures adopted to the weight data.
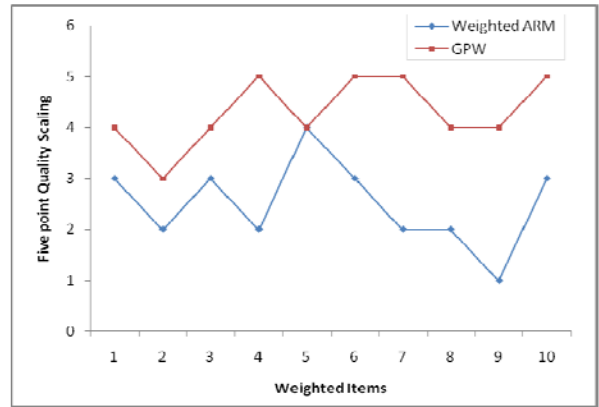


Fig 11: Five-point quality scaling

It is inferred from the above table that GPW has got highest quality scaling grade for most data. The result exhibits GPW is powerful and has more meaningful result than traditional weighted association rule mining. It is also understood from the GPW measurement that frequency of computation will be less than weighted ARM. Due to high frequency items, changes in criteria and changes in profit for item situation only needed to recompute the GPW, else can reuse the data with previous weight measures.
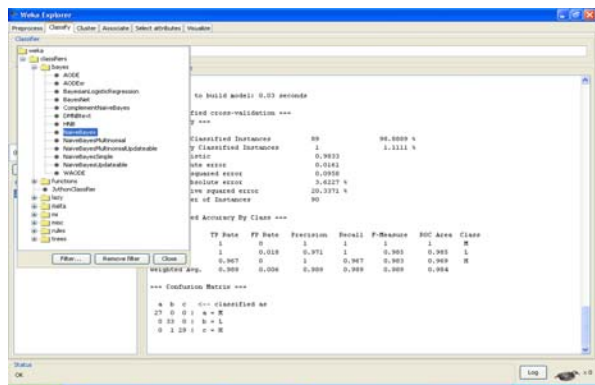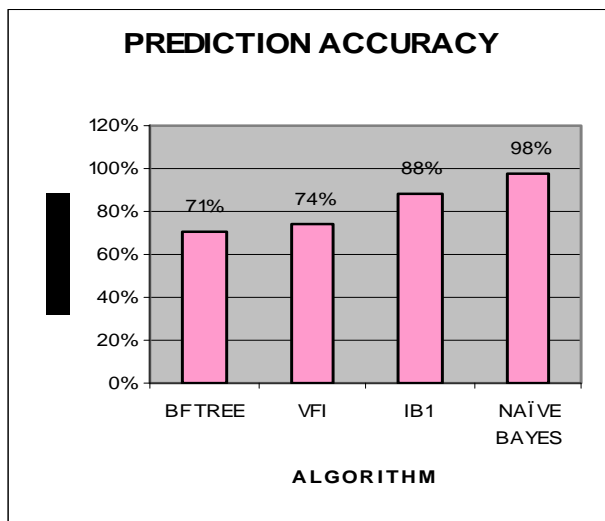
I    Training And Testing

Fig 12: Selection of Naïve Bayes from the Classifier Tab

The comparative evaluation results are summarized in Table 2. The performance of the three models were evaluated based on the three criteria, the prediction accuracy, learning time and error rate are illustrated in Figures 13, 14 and 15.

| Evaluation Criteria | NAÏVE BAYES | IB1 | VFI | BF TREE |
|---|---|---|---|---|
| Timing to build model (in secs) | **0.02 sec** | 0.03 sec | 0.5 sec | 0.36 sec |
| Correctly classified Instances | **98** | 88 | 74 | 71 |
| Incorrectly Classified Instances | **2** | 10 | 23 | 26 |
| Prediction Accuracy (%) | **98.8%** | 88.8% | 74.4% | 71.1% |

Table 2: Comparative results of the classifiers



PREDICTION ACCURACY

Fig 13: Prediction Accuracy

As shown in Figure 13, Naïve Bayes predicts better than other algorithms. Among the four classifiers used for the experiment, the IB1 classifier is more or less the same prediction accuracy. The accuracy rate of BF Tree and VFI classifier is the lowest among the four machine learning techniques.
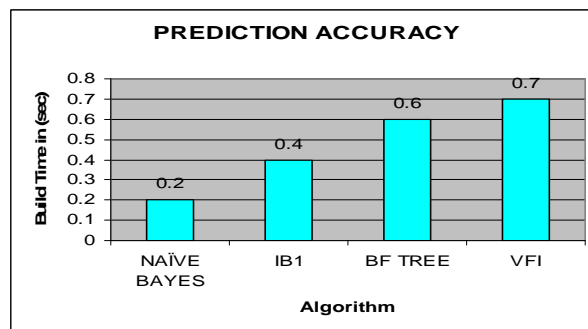


Fig 14: Time Taken to Build

As shown in Figure 14, the build times of the four schemes are under consideration. The Naïve Bayes, the probabilistic classifier tends to learn more hastily for the given dataset. There is a little statistical variation in the time taken to build the IB1, VFI and BF Tree classifier model and probabilistic model.
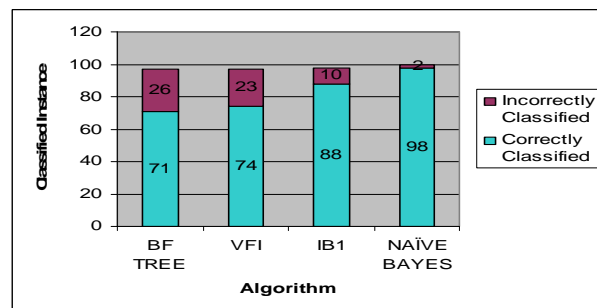


Fig 15: Classified Instance

Figure 15 shows the correctly classified instances and incorrectly classified instances from number of instances of four classifiers.

VI. CONCLUSION

The objective of study is to implement the global profit weight measure and test the performance of the algorithm with traditional weighted association rule mining. The implementation is expressed as step wise visual presentation and its performance is measured with weighted ARM.

The accuracy is classified using classification algorithm such as Navie Bayes, VFI, BF Tree and IB1and the results are compared using WEKA. It can be concluded from the study result that GPW is required high computation power to generate the weight. The result is compromised with its quality. According to the research

problem, the calculated weight can be reused it for many times and as required.

## REFERENCES

[1] C. Aggarwal. Towards long pattern generation in dense databases. SIGKDD Explorations, 3(1):20{26, 2001.

[2] .R. Agrawal, C. Aggarwal, and V. Prasad. Depth First Generation of Long Patterns. In 7th Int'l Conference on Knowledge Discovery and Data Mining, Aug. 2000.

[3] .R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in largedatabases.In Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data, pages 207-216. ACM Press, 1993.

[4] Y. Bastide, R. Taouil, N. Pasquier, G. Stumme, and L. Lakhal. Mining frequent patterns with counting inference. SIGKDD Explorations, 2(2), Dec. 2000.

[5] R. J. Bayardo. Efficiently mining long patterns from databases. In ACM SIGMOD Conf. Management of Data, June 1998.

[6] Barber and H.J. Hamilton. Extracting share frequent itemsets with infrequent subsets. Data Mining and Knowledge Discovery, 7(2):153-185,2003

[7] F. Bodon. A fast apriori implementation. In Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations, 2003.

[8] B. Goethals. E±cient Frequent Pattern Mining. PhD thesis, transnational University of Limburg, Belgium, 2002.

[9] K. Gouda and M. J. Zaki. Efficiently mining maximal frequent itemsets. In 1st IEEE Int'l Conf. on Data Mining,

[10] Nov. 2001.

[11] G. Grahne and J. Zhu. High performance mining of maximal frequent itemsets. In 6th International Work- shop on High Performance Data Mining, May 200

## AUTHORS PROFILE

ASHA RAJKUMAR M.phil, Reasearch Scholar.

Computer Science Department
P.S.G.R.Krishnammal College for Women, Coimbatore, India

G.SOPHIA REENA M.C.A., Mphil.,

HOD, BCA Department
P.S.G.R.Krishnammal College for Women
Coimbatore, India