# A systematic review on Mining the infrequent item sets from frequent patterns based on FTP (Fast Progress Technique)

**\*T.Sudhan; \*\*C.Anuradha&\*\*\*Dr.C.Nalini**

\*M.Tech-Student, Dept of CSE, Bharath University, Email: suviandvisu@gmail.com
\*\*Assistant Professor, Dept of CSE, Bharath University,
\*\*\*Professor, Dept of CSE, Bharath University,

## Abstract

*Data mining, which is the exploration of knowledge from the large set of data, generated as a result of the various data processing activities. Pattern Mining is a very important task in data mining. Whereas, frequent pattern mining has been a focused theme in data mining research for over two decades. In past literature has been dedicated to this research andtremendous progress has been made especially mining the infrequent item sets from frequent patterns based on FTP (fast progress technique), as well as their broad applications. In this paper we present the reviews on different pattern growth methods from infrequent item sets to frequent pattern based FTP techniques.*

**Keywords:**
Data Mining, Frequent, FTP, Infrequent, Itemsets and Pattern.

## Introduction

Frequent weighted itemsets represent correlations frequently holding in data in which items may weight differently. However, in some contexts, e.g., when the need is to minimize a certain cost function, discovering rare data correlations is more interesting than mining frequent ones. **Luca Cagliero and Paolo Garza et.al. 2014**, tackled the issue of discovering rare and weighted itemsets, i.e., the infrequent weighted itemset (IWI) mining problem. Two novel quality measures were proposed by them to drive the IWI mining process. Furthermore, two algorithms that perform IWI and Minimal IWI mining

efficiently, driven by the proposed measures, are presented. Experimental results shows efficiency and effectiveness of the proposed approach [1].

In particular, they focused on attention for two different IWI-support measures: (i) The IWI-support-min measure, which relies on a minimum cost function, i.e., the occurrence of an itemset in a given transaction is weighted by the weight of its least interesting item, (ii) The IWI-support-max measure, which relies on a maximum cost function, i.e., the occurrence of an itemset in a given transaction is weighted by the weight of the most interesting item. They noticed that, when dealing with optimization problems, minimum and maximum are the most commonly used cost functions. Hence, they are deemed suitable for driving the selection of a worthwhile subset of infrequent weighted data correlations. Specifically, the following problems have been addressed [1, 2, and 3]:

A. IWI and Minimal IWI mining driven by a maximum IWI-support-min threshold, and

B. IWI and Minimal IWI mining driven by a maximum IWI-support-max threshold.

Since this paper reviews about the mining infrequent item sets from frequent patterns based on FTP technique, here the content 2 explains few literature survey based on the patterns designed and explained so far and then the content 3 reviews on experimental views which relies on infrequent item sets from frequent patterns.

## 2. Literature Survey

One of the results of the literature survey shown here explains, the work in the designing of tree structures for use in association rule mining. Frequent Pattern Tree (known as FP-Tree) was first introduced by **Han et al (2000)** is based upon a tree representation of frequent itemsets. It compresses a large database of transaction into a compact, Frequent Pattern Tree (FP-Tree) structure thereby eliminating candidate generation. It gets its name from the fact that only frequent itemsets are used to build the tree. It requires just two scans of the data database, it avoids the candidate itemset generation, and it allows repeated mining of the same data without further database scan.

**Han et al (2004)** represented a novel frequent-pattern tree (FP-Tree) structure, which is an extended prefix-tree structure for storing compressed, vital information about frequent patterns and develop an efficient FP-Treebased mining method, FP-growth, for mining the complete set of frequent patterns by pattern fragment growth. This method avoids the computational intensive process of candidate generation and testing. This approach substantially reduces search time.

**Cheung et al (2004)** proposed two new algorithms namely to show that their methods do better than the previously proposed Itemset-Loop algorithm. They also proposed the mining of "N-most interesting k-itemsets with item constraints". **Burdick et al (2005)** presented a new algorithm for mining maximal frequent itemsets from a transaction database. The search strategy of the algorithm integrates a depth-first traversal of the itemset lattice with effective pruning mechanisms that significantly improve mining performance.

**Zaki et al (2005)** presented an efficient algorithm for mining all frequent mining closed itemsets. It enumerates closed sets using a dual itemset-tidset search tree, using an efficient hybrid search that skips many levels.**Grahne et al (2005)** presented a novel FP-array technique that greatly reduces the need to traverse FP-Trees, thus obtaining significantly improved performance for FP-Tree-based algorithms. Their technique works especially well for sparse data sets. Furthermore, they presented new algorithms for mining all, maximal, and closed frequent itemsets. Their algorithm uses the FP-Tree data structure in combination with the FP-array technique efficiently and incorporates various optimization techniques.

**Bonchi et al (2005)** defined a data mining query language, supported by a system that can optimize constraint-based data mining queries. They have invented ExAnte, a simple yet effective preprocessing technique for frequent pattern mining.**Song et al (2006)** presented a novel algorithm for mining complete frequent itemsets. This algorithm is referred to as the Transaction Mapping (TM) algorithm. In this algorithm, transaction IDs of each itemset are mapped and compressed to continuous transaction intervals in a different space and the counting of itemsets is performed by intersecting these interval lists in a depth-first order along the lexicographic tree.

**Wang Jianyong et al (2007)** proposed an efficient algorithm, for mining frequent closed sequences without candidate maintenance. It adopts a novel sequence

closure checking scheme called BI-Directional Extension (BIDE) and prunes the search space more deeply compared to the previous algorithms by using the BackScan pruningmethod. An efficient algorithm named Apriori -growth based on Apriori algorithm and the FP-Tree structure is presented by **Wu Bo et al (2008)** to mine frequent patterns. The advantage of the Apriori -growth algorithm is that it doesn'tneed to generate conditional pattern bases and sub-conditional pattern treerecursively. Computational results show the Apriori -growth algorithmperforms faster than the Apriori algorithm, and it is almost as fast as FPgrowth,but it needs smaller memory.

**Gupta et al (2008)** proposed a comprehensive evaluation framework to compare different approximate frequent pattern mining algorithms. The key idea is to select the optimal parameters for each algorithm on a given data set and use the itemsetsgenerated with these optimal parameters in order to compare different algorithms.Mining frequent patterns from XML documents can be recast as mining frequent tree structures from a database of XML documents.

**Tan et al(2008)**modeled a database of XML documents as a database of rooted labelled ordered subtrees. **Li Yu-Chiang et al (2008)** introduced the Isolated ItemsDiscarding Strategy (IIDS), which can be applied to any existing level-wiseutility mining method to reduce candidates and to improve performance. **JinRuoming et al (2008)** proposed a set of novel regression-based approaches toeffectively and efficiently summarize frequent itemset pattern. Specifically,they have shown that the problem of minimizing the restoration error for a setof itemsets based on a probabilistic model corresponds to a nonlinearregression problem. Mining frequent patterns from XML documents can berecast as mining frequent tree structures from a database of XML documents.

**Huang et al (2008)** discussed the problem on mining frequentepisodes in a complex sequence. They extended the previous algorithm forepisode mining from complex sequences. Furthermore, a memory-anchoredalgorithm is introduced to the mining task. **Xu Yu Jeffrey et al (2008**)presented a data mining proxy approach that reduces the I/O costs to constructan initial tree by utilizing the trees that have already been resident in memory.The tree they constructed is the smallest for a given data mining query.**Yun Unil (2008)** proposed weighted frequent pattern mining withlength decreasing support constraints. Their main approach is to push weightconstraints and length decreasing support constraints into the pattern growthalgorithm. For pruning, they proposed the notion of the Weighted SmallestValid Extension (WSVE) property with/without Minimum Weight.

**Fan et al(2008)** proposed a new and different method to build a decision tree thatpartitions the data onto different nodes. Then at each node, it directlydiscovers a discriminative pattern to further divide its examples into purersubsets. **Lee et al (2008)** proposed an efficient algorithm, called IntertransactionClosed patterns Miner (ICMiner), for mining closed intertransactionitemsets. Their proposed algorithm consists of two phases. First,scan the database once to find the frequent items. For each frequent itemfound, the ICMiner converts the original transaction database into a set ofdomain attributes, called a data set. Then, it enumerates closed

intertransactionitemsets using an itemset–data set tree, called an ID-tree. **Raissiet al (2008)** aimed at extending the non-derivable condensed representation infrequent itemset mining to sequential pattern mining. They started by showinga negative example: in the context of frequent sequences.

**Lin Chun-Wei et al (2009)** attempted to modify the Fast UpdatedFP-Tree (FUFP-tree) construction based on the concept of pre-large itemsets.Pre-large itemsets are defined by a lower support threshold and an uppersupport threshold. It does not need to rescan the original database until anumber of new transactions have been inserted. Woo Ho**Jin et al (2009)**proposed a method of tracing the set of MFIs instantly over an online datastream. The method, namely estMax, maintains the set of frequent itemsets bya prefix tree and extracts all MFIs without any additional superset/subsetchecking mechanism. **Gomez et al (2009)** proposed a much powerfullanguage, based on regular expressions, where the basic elements areconstraints defined over the (temporal and non-temporal) attributes of theitems to be mined. **Zailani et al (2010)** proposed a scalable tire-basedalgorithm; this algorithm generates the significant patterns using intervalsupport and determines its correlation. Experiments with the real data setsshow that this algorithm can discover highly positively correlated andsignificant of least association.

**Chen et al (2010)** presented a three-strategy adaptive algorithm,Bitmap Itemset Support Counting (BISC), is presented. The core strategy,BISC1, is used in the innermost steps of the recursion. For a database $D$ withonly $s$ frequent items, a depth-first approach need up to $s$ levels of recursionsto detect all the

FIs. **Abdullah et al (2010)** proposed a scalablemodel called Critical Least Association Rule (CLAR) to discover thesignificant and critical least association rules. Experiments with a real andUCI data sets show that the CLAR can generate the critical least associationrules, up to 1.5 times faster and less 100% complexity than benchmarkedFP-Growth. **Jose et al (2010)** focuses mostly on the case were labels of thenodes are nonexistent or unreliable, and discuss algorithms for closure-basedmining that only rely on the root of the tree and the link structure. Theyprovided a notion of intersection that leads to a deeper understanding of thenotion of support-based closure, in terms of an actual closure operator.

**Lucchese et al (2010)** formalized the problem of discovering theTop-K patterns from binary data sets in the presence of noise, as theminimization of a novel cost function. According to the MinimumDescription Length principle, the proposed cost function favors succinctpattern sets that may approximately describe the input data. **Kiran et al (2010)** exploited the notion of "item-to-pattern difference" and propose multiple min_sup based FP-growth-like approach to efficiently discover rare association rules. **Tseng et al (2010)** presented an efficient algorithm, namely Utility Pattern Growth (UP-Growth), for mining high utility item sets with aset of techniques for pruning candidate item sets. The information of highutility item sets is maintained in a special data structure named Utility PatternTree (UP-Tree) such that the candidate itemsets can be generated efficientlywith only two scans of the database.

**Vo Bay et al (2011)** proposed a new method with a supportingstructure called Attributes Itemset Object identifications - tree (AIO-tree) formining FIs from multidimensional databases. This method need not transformthe database into the transaction database, and it is based on the intersectionsof object identifications for fast computing the support of itemsets. **Floratou etal (2011)** presented a new algorithm called Flexible and Accurate MotifDetector (FLAME). FLAME is a flexible suffix-tree-based algorithm that canbe used to find frequent patterns with a variety of definitions of motif(pattern) models.

In this literature survey, we have gone through a lot about frequent patterns, but the main systematic theme of this review paper to explain about the experimental mining in infrequent items sets from frequent patterns base on FTP technique is as follows in experimental observation.

# 3. Experimental Observations:

Mining infrequent patterns is a challenging endeavour because there are an enormous number of such patterns that can be derived from a known data set. More exclusively, the key issues in mining infrequent patterns are: (1) how to identify interesting infrequent patterns, and (2) how to efficiently discover them in large data sets. To get a different perspective on various types of interesting infrequent patterns, two connected conceptions are negative patterns and negatively correlated patterns. The negative itemsets and negative association rules are collectively known as negative patterns. Infrequent patterns, negative patterns, and negatively correlated patterns are three closely related concepts. Although infrequent patterns and negatively correlated patterns refer only to itemsets or rules that contain positive items, while negative patterns refer to itemsets or

rules that contain both positive and negative items [3, 4 and 5]. Here, will correlate the major pattern growth for infrequent and frequent mining techniques by FPT.

## 3.1 Pattern growth methods for Infrequent and frequent pattern mining

Infrequentitemsets are given by all itemsets that are not extracted by standard frequent itemset generations algorithms such as Apriori and FP-growth. Since the number of infrequent patterns exponentially large, especially for sparse, high dimensional data, techniques developed for mining infrequent patterns focuson finding only interesting infrequent patterns.[5,6].

*Mining Negative Patterns:* Transaction data can be binarized by augmenting it with negative items. By applying existing frequent itemset generation algorithm such as Apriori on the augmented transactions, all the negative itemsets can be derived. Such an approach is feasible only if a few variables are treated as symmetric binary.

*Support Expectation:* Another class of techniques considers an infrequent pattern to be interesting only if its actual support isconsiderably smaller than its expected support. For negatively correlated patterns, the expected support iscomputed based on the statistical independence assumption. Two alternative approaches for determiningthe expected support of a pattern using (1) a concept hierarchy and (2) a neighborhood-based approach known as indirect association.

*Support Expectation Based on Concept Hierarchy:*Objective measures alone may not be sufficient to eliminate uninteresting infrequent patterns. For example,support bread and laptop computer are frequent items. Even though the itemset {bread, laptop computer} is infrequent and perhaps negatively correlated, it is not interesting because their lack of support seems obvious to domain experts. Therefore, a subjective approach for determining expected support is needed to avoid generating such infrequent

patterns. In the preceding example, bread and laptop computers belong to two completely different product categories, which is why it is not surprising to find that their support is low.

***Support Expectation Based on Indirect Association:***Consider a pair of items, (a, b), that are rarely bought together by customers. If a and b are unrelated items such as bread and DVD player, then their support is expected to be low. On other hand, if a and b are related items, then their support is expected to be high. The expected support was previously computed using a concept hierarchy. Here, an approach for determining the expected support between a pair of items by looking at other items commonly purchased together with these two items. Indirect association has many potential applications. In the market basket domain, a and b may refer to computing items such as desktop and laptop computers. In text mining, indirect association can be used to identify synonyms, antonyms, or words that are used in different contexts. For example, given a collections of documents, the word data may be indirectly associated with gold via the mediator mining. This pattern suggests that the word mining can be used in two different contexts − data mining versus gold mining.

### 3.1.1 Weighted Frequent Itemsets Mining

Researchers have proposed weighted frequent itemset mining algorithms that reflect the significance of items.The foremost focus of weighted frequent itemset mining isconcerns satisfying the downward closure belongings.Every weighted association rules mining algorithmssuggested so far have been based on the Apriori algorithm.Nevertheless, pattern growth algorithms are much moreefficient than Apriori based algorithms. An efficientweighted frequent itemset mining algorithm is the mainapproach used to push weight constraints into the patterngrowth algorithm and provide ways to keep the downwardclosure assets. WFIM accepts a rising weight orderedprefix tree. The tree is

traversed bottom-up because theprevious matching cannot maintain the downward closureproperty. A support of each itemset is usually decreased asthe length of an itemset is enlarged, but the weight has aunusual characteristic. An itemset which has a low weightsometimes can get a higher weight after adding anotheritem with a higher weight, so it is not guaranteed to keepthe downward closure property [7].

### 3.2 FTP preserving in Infrequent and frequent pattern mining in large database

### 3.2.1 Frequent pattern mining in large database
If the frequent-item projections of transactions in the database can be held in the main memory, then they can be organized as shown in ***Fig. 1***. All items in frequent-item projections are sorted according to the *F-list*. For example, the frequent-item projection of transaction is listed as *cdeg*. Every occurrence of a frequent item is stored in an entry with two fields: an *item-id* and a *hyper-link*. A*headertableH*is created, with each frequent item entry having three fields: an *item-id*, a *support count*, and a *hyperlink*. When the frequent-item projections are loaded into the memory, those with the same first item (in the order of the *F-list*) are linked together by the hyper-links into a queue, and the entries in header table *H* act as the heads of the queues. For example, the entry of item *a* in the header table *H* is the head of the *a*-queue, which links frequent-item If the frequent-item projections of transactions in the database can be held in the main memory, then they can be organized as shown in Fig. 1. All items in frequent-item projections are sorted according to the *F-list*. For example, the frequent-item projection of transaction 100 is listed as *cdeg*. Every occurrence of a frequent item is stored in an entry with two fields: an *item-id* and a *hyper-link*. A *header table H* is created, with each frequent item entry having three fields: an *item-id*, a *support count*, and a *hyperlink*. When the frequent-item projections are loaded into the memory, those with the same first item (in the order of the *F-list*) are linked

together by the hyper-links into a queue,nd the entries in header table *H* act as the heads of the queues. For example, the entry of item *a* in the header table *H* is the head of the *a*-queue, which links frequent-item.
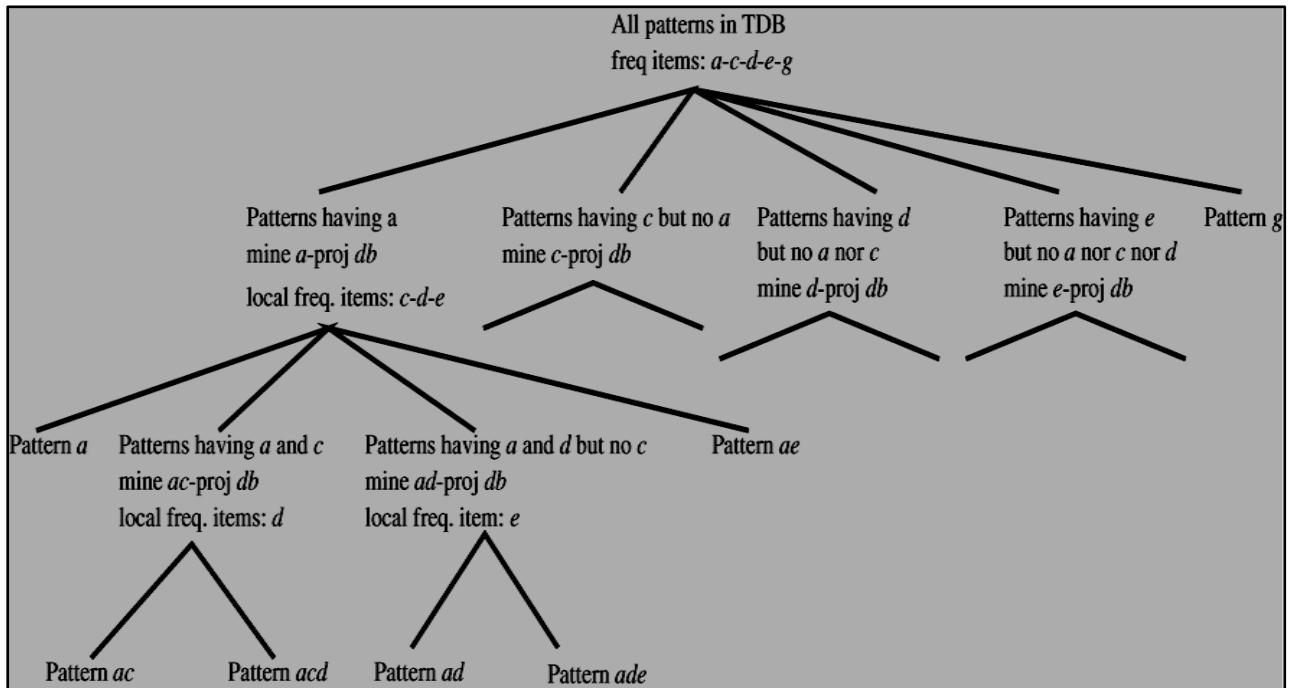


*Fig 1: Divide-and-conquer frequent patterns.*

### 3.2.2 Infrequent pattern mining in large database

Even though infrequent pattern mining is still an emerging research field and has been studied for a decade, there are still many unsolved topics that can be explored, such as how to further control the number of generated candidate and how to improve the efficiency of the mining process by providing more targeted candidates, etc [8]. It is not hard to generate a large number of infrequent sub trees, but the challenge comes from how to determine which infrequent sub trees are most interesting and valuable to end users and their applications [9].
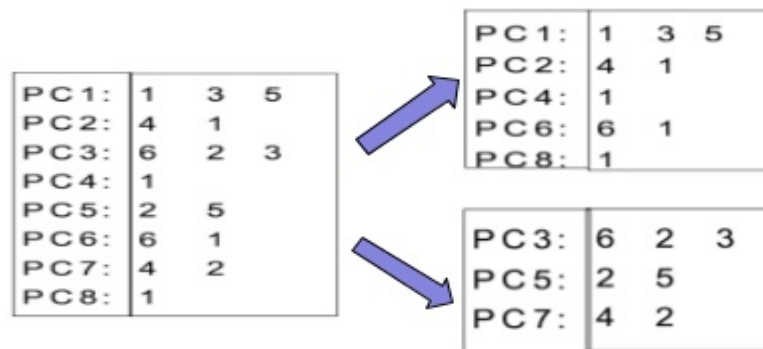


*Fig 2: Divide-and-conquer of infrequent patterns.*

To find minimal infrequent item [10] sets developed for finding minimal unique item sets. Infrequent item setwhich has an infrequent proper subset is redundant, since the former can be deduced from the latter. i) computing the support of each item, which is needed to produce a rank-ordering of the viable items by support ii) determining the viability of each item pruning some of the items from consideration iii) computing the support set of each viable item resulting in a memory efficient representation of the lists of TIDs for each support set IWI Miner is a FP-growth-like mining algorithm that performs projection-based item set mining.

FP-growth mining steps:
1. FP-tree creation and (b)
2. Recursive item set mining from the FP tree index.
3. IWI Miner discovers infrequent weighted item sets instead of frequent (unweight) ones.
Modifications with respect to FP-growth have been introduced:
(i) A novel pruning strategy for pruning part of the search space early and
(ii) a slightly modified FP tree structure, which allows storing the IWI-support value associated with each node.

**Infrequent weighted item set Algorithm**
Input- weighted transaction dataset and support value)
**IWI (T, E)**
1) F=0
2) Count item IWI (T)
3) Construct FP tree
4) For all weighted transaction
5) Calculate Equivalent transaction
6) For all transaction create and insert into FP tree
Output – Set of satisfying E

MIWI Miner focuses on generating only minimal infrequent patterns, the recursive extraction in the MIWI Mining procedure is stopped as soon as an infrequent item set occurs. It finds both the infrequent item sets and minimal infrequent item set mining.

**IWI mining (T, E, P)**
1) F=0 initialization
2) Create header table holds for all items i in tree
3) Generate a new item set I with prefix and support of item i
4) I – Infrequent item
5) Construct I as conditional pattern and FP tree
6) Select the infrequent items from the set
7) Remove from Tree and finally apply recursive mining

## 4. Practical Applications

Infrequent patterns can be used in many applications.

- In text mining, indirect associations can used to find synonyms, antonym or words that are used in different contexts. For example, the word *data* might be indirectly associated with the word *gold*, using the mediator *mining*.
- In the market basket domain, indirect associations can be used to find competing items, such as *desktop computers* and *laptops*, which states that people whom buys desktop computers won't buy laptops.
- Infrequent patterns can be used to detect errors. For example, if *{Fire = Yes}* is frequent, but *{Fire = Yes, Alarm = On}* is infrequent, then the alarm system probably is faulting [11].

- When evaluating Weka [12], we could not find any signs of an infrequent pattern classifier.
- Searching for outliers in data stream is an important area of research in the world of data mining with numerous applications, including credit card fraud detection, discovery of criminal activities in electronic commerce, weather prediction, marketing and customer segmentation [13].The outliers can be detected from infrequent patterns.
- Intrusion detection in wireless networks has become a vital part in wireless network security systems with wide spread use of Wireless Local Area Networks (WLAN). Intrusion detection can be done by the infrequent patterns [14].

## 5. Results and Discussions

Infrequent pattern mining is still an emerging research field and has been studied for a decade, there are still many unsolved topics that can be analysed, such as how to further control the number of generated entity and how to improve the abstraction of the mining process by providing more targeted entity, etc. It is not hard to develop a large number of infrequent sub trees, but the challenge comes from how to determine which infrequent sub trees are most interesting and valuable to end users and their applications. By reading this review paper one can get a basic knowledge on mining the infrequent and frequent patterns based on FTP and this work acts as a basis for the future work that has been done on this emerging area.

## 6. Conclusion

This review paper mainly concentrated on mining the infrequent item sets from frequent patterns based on FTP techniques. To get a different perspective on various types of interesting frequent and infrequent patterns, related and correlated patterns for large database is mentioned above. It explains what are patterns, item set, and association rule. This paper makes review on different papers related to frequent and infrequent patterns and rare item sets and also gives the knowledge on FTP proposed for mining infrequent item sets from frequent patterns. This also explainsabout different application areas where these infrequent patterns are used. This research work lucubrate the algorithms for mining infrequent patterns from frequent item sets and which becomes the basis for the future work going to be done in this area.

## 7. Future Work

Mining infrequent patterns is a challenging attempt, because there is an extent number of such patterns that can be derived from a given data set. More specifically, the key issues in mining infrequent item sets from frequent patterns are:
(1) How to identify infrequent patterns,
(2) How to identify the interestingness of those patterns
(3) How to efficiently discover them in large data sets.

Thus to get a different perspective on various types of interesting mining the infrequent item sets from frequent patterns, are related concepts done by FTP is an enormous process in this emerging field.

**References**

[1.] BalazesRacz," nonordfp: An FP-Growth Variation without Rebuilding the FP-Tree**",** 2nd Int'l Workshop on Frequent ItemsetMiningImplementations FIMI2004J.

[2.] Conference on Intelligent Computation Technology and Automation, pp. 394- 397, 2012.

[3.] G. Cong, A.K.H. Tung, X. Xu, F. Pan, and J. Yang, "Farmer: Find- ing Interesting Rule Groups in Microarray Datasets," Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '04), 2004.

[4.] Han, J. Pei, and Y. Yin, "Mining Frequent Patterns without Candidate Generation," Proc. ACM SIGMOD Int'l Conf. Management of Data,pp. 1-12, 2000.

[5.] International Journal of Artificial Intelligence & Applications (IJAIA), Vol.1, No.4, pp. 133 – 143, October 2010.

[6.] JiaRong Bit (hons) Advanced Pattern Mining for Complex Data Analysis submitted in fulfillment of the requirements for the degree of Doctor of Philosophy, Deakin University, August 2012.

[7.] Luca Cagliero and Paolo Garza "Infrequent Weighted Itemset Mining using Frequent Pattern Growth", IEEE Transactions on Knowledge and Data Engineering, pp. 1- 14, 2013.

[8.] M.Hamsathvani1, D.Rajeswari2, R.Kalaiselvi 3Survey on Infrequent Weighted Itemset Mining Using FP Growth, 10.15662/ijareeie.2014.0311027.

[9.] M.L. Antonie, O.R. Zaiane, and A. Coman, "Application of Data Mining Techniques for Medical Image Classification," Proc. Second Intl. Workshop Multimedia Data Mining in Conjunction with seventh ACM SIGKDD (MDM/KDD '01), 2001.

[10.] R. Agrawal, T. Imielinski, and Swami, "Mining Association Rules between Sets of Items in Large Databases," Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '93), pp. 207-216, 1993.

[11.] ShipraKhare et al, / (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 5 (2) , 2014, 1642-1647

[12.] Soumadip Ghosh, Sushanta Biswas, Debasree Sarkar, ParthaPratim Sarkar "Mining Frequent Itemsets Using Genetic Algorithm",

[13.] W. Wang, J. Yang, and P.S. Yu, "Efficient Mining of Weighted Association Rules (WAR)," Proc. Sixth ACM SIGKDD Int'l Conf. Knowledge Discovery and data Mining (KDD '00), pp. 270-274, 2000.

[14.] Xin Li, Xuefeng Zheng, Jingchun Li, Shaojie Wang "Frequent Itemsets Mining in Network Traffic Data", 2012 Fifth International.