

Review on Frequent Itemset Mining Via Transaction Splitting

Miss.Pooja Purohit¹, Prof.Sonal Patil²

¹Computer Engineering Department,G.H. Rasoni College of Engineering,
Jalgaon,Maharashtra,India

² Computer Engineering Department,G.H. Rasoni College of Engineering,
Jalgaon,Maharashtra,India

Abstract: Nowadays, businesses are evolving. For having business people needs to deal with much amount of data and this data needs to be delicate and confidential. So, to secure and preserve our data there are plenty of technologies used one of them is Data Mining. Data Mining is the technique in which it tries to find out interesting patterns or knowledge from database such as association or correlation etc. Frequent Itemset Mining is the critical problem in data mining. The frequent can contains valuable and research purpose. Frequent itemsets are items or patterns like itemset, substructures or subsequences that occurs frequently in transaction. To find out frequent itemset there are many Frequent Itemset Mining Algorithms used such as Apriori, FP-growth, Elcat. The famous two most important algorithm, to find Frequent itemset are Apriori and FP-growth. Apriori is a candidate set generation and-test algorithm. It needs multiple database scans.It have two steps: It find all itemsets which have minimum support frequent itemsets(candidate list) and generate frequent itemsets list . FP-growth which does not generates candidate set. Compared with Apriori, FP-growth is less time consuming algorithm. But it enforcing the limit by truncating transactions, if a transaction has more items than the given limit, then deleting items until its length is under the limit. FP-growth is more faster than apriori algorithm. But this all only consider frequent itemset from large transaction. It doesn't consider the useful or high utility itemset from large transaction. For privacy preservation it uses the differentially privacy method.

Keywords - Frequent Item Mining, Apriori algorithm, FP- growth

I. Introduction

Data Mining is the method which find out the hidden data in a database. It is also called as data determination or data analysis and deductive finding out. There is a great deal of knowledge to extract the associations from data. To identify the correlated set of item in database association rules are used. Data Mining uses different kinds of techniques which are combined from database technologies and many more.

Market basket analysis is the one of the most famous example of association rule mining. In this, market analysts focused on discovering frequently purchased items by the consumer. So it is easy for arrangement of items according to their sales by organization. Association rule mining is the method of finding interesting or similar relations between the variables or items in large transaction of database. Association rule is using two measures support and confidence to identify the most important relationships. Support is an sign of how frequently the itemset appear in the database and it is also set of preconditions.

Confidence is an mark of finding how frequently the rule has been found to be true. It uses minimum support and confidence which are user defined. Association rule used in many application areas including web usage mining,intrusion detection , continuous production and bioinformatics.

Discovering useful patterns hidden in database plays an critical role in different data mining jobs such as frequent pattern mining , high utility pattern mining .Among the all ,frequent pattern mining is a basic research topic that has been used to different database having long transactions. It is used in the analysis of customer transactions in retail research where it is marked as market basket analysis and also used to identify purchase items of the consumer. Given a database, in which each of transaction has a set of items where Frequent Itemset Mining used to find out itemsets that occurs in transactions more than a given user specified threshold.

The data is perceptive (e.g web browsing history and medical records of patients) the frequent itemsets detection can provide. Releasing that detected frequent itemsets can cause threats to individual privacy. But limitations of frequent itemset mining are that only consider frequently occurring items in a transaction database above the user specified frequency threshold, without considering the quantity and profit of items. The quantity and utility are important for real world decision problem.

This paper marks the frequent and weighted itemsets discovery. Most of the methods in finding frequent itemsets which designed for traditional databases they are apriori and FP-growth algorithm.

II. Related Work

Many of researchers have been proposed to solve the privacy preserving FIM problem from different ways. Mining the frequent patterns is mine in many different kinds of databases such as transaction database, time-series databases, and many other. These databases have been investigated in data mining research. Many of the previous examination accept an *Apriori* as like candidate set which is generation-and-test method. The candidate set generation is costly. In this study, the FP-growth, it structure an extended prefix-tree structure which is used for storing compressed and important information about frequent patterns.

Main purpose is that the resulted frequent itemsets itself does not leak private information and achieve differential privacy. The k -anonymity model for protecting privacy in [2] and in [12] which propose an algorithm to publish anonymised frequent itemset. Both studies don't satisfy differential privacy as well as they cannot provide sufficient privacy protection from attackers having background knowledge. Attribute [3] introduced l -diversity, a framework which gives stronger privacy guarantees and shows the weak points of k -anonymity.

In [4] proposed Apriori & AprioriHybrid algorithms which are fast algorithm for mining association rule. These both of compared with previous algorithms and gives excellent performance for large database with transactions, but it generates candidate set which is costly to handle these which is costly to handle these.

In [5] Introduces FP growth algorithm, which is mining frequent pattern without candidate generation and as we seen in [4] apriori algorithm performs mining fastly with candidate set generation, which is costly. In [5] FP-growth it uses FP-tree as a data structure to store large database which compressed in small data structure. Here, in [5] the algorithm is scalable and efficient than apriori algorithm.

In this [7], algorithm solves the frequent item set mining problem that they find all item set whose support exceeds a threshold. The benefit of this algorithm is that it achieves better F-score unless k is small. C. Zeng, J. F. Naughton, and J.-Y. Cai,[7] proposed an apriori algorithm for large transaction. Major problem is for long transaction which contains many items. Its truncating the long transactions means it limiting the transaction. Deleting the items until the transaction is under the limit if transaction has more than a specified number of items. It must be done in a differentially private way. It also discarding items from transactions gives a new source of error. In transaction truncating, the more frequent subsets are kept and other items which are not frequent are truncated.

In this [8] it solve the problem of association rule mining algorithm which gives a privacy preserving scalar product protocol as well as gives an efficient protocol for computing scalar product which preserve privacy of the individual transactions. In [9] Clifton and Kantarcioglu, which addresses the problem as a secure multi-party computation and consider the database as a horizontally partitioned.

L.Bonomi [6] proposed Frequent sequential pattern mining. It is a central task in many fields as like biology and finance. In this paper, it provides the provable and formal guarantees of privacy by studying the sequential pattern mining problem under the differential privacy framework. In this paper, proposed a two-phase algorithm. This algorithm mining both prefixes and substring patterns. In the first phase, construct a model-based prefix tree. This is used to mine a candidate set of substring patterns and its prefixes. In next phase, substring patterns is refined. Here it transformed the original data to reduce the perturbation noise.

W.K.Wong [10] proposed Outsourcing association rule mining to an outside service provider which gives many important benefits to the data owner. These include (i) release from the high mining cost, (ii) minimize the demands in resources, and (iii) for multiple distributed owners effectively centralized the mining and also provide security. In this paper, it develop an effective and efficient encryption algorithm and performs a single pass over the database. It is applicable for the application which sends the streams of transactions to the service provider.

In [11] it present the set of randomization operators to limit privacy which margin the FIM. Proposed new algorithm which discovers the frequent patterns in sensitive data and adopted two mechanism techniques i.e. exponential & Laplace noise-addition mechanism. These are efficient in context of frequent item mining.

[12] Proposes algorithm Privbasis which perform frequent itemset mining with differential privacy with the help of minimum support threshold. An item set that found in transaction is frequently than minimum support threshold and subset of some basis with differential privacy guarantee.

III. Methodology

Apriori is a breadth first search algorithm. It generates the candidate set. So, it also called as generation-and-test algorithm. It needs only one database scans if the maximal length of frequent itemsets is l . It Support count is expensive, due to generation of candidate set and requires multiple database scans (I/O). It consist of two steps: Find all itemsets which have minimum support frequent itemsets. Use that frequent itemsets to generate rules. The Apriori property follows a two step process:

- Join step: C_k is generated by joining L_{k-1} with itself.

- Prune step: Any (k-1)-itemset that is not frequent cannot be a subset of a frequent k-itemset

1.1 The Apriori Algorithm

C_k : Candidate itemset of size k
 L_k : frequent itemset of size k
 $L_1 = \{\text{frequent items}\};$
 For (k= 1; $L_k \neq \emptyset$; k++) do begin
 C_{k+1} = candidates generated from L_k ;
 for each transaction t in database do
 increment the count of all candidates in C_{k+1} that are contained in t .
 L_{k+1} = candidates in C_{k+1} with min_support
 End
 Return $\cup_k L_k$

1.1.1 Limitation

- Needs so many iterations of the data and uses uniform minimum support threshold. Problem in finding rarely occurring events
- Due to generation of large number of candidates it requires large memory space. Execution is more as time is wasted in producing candidate every time
- Execution time is more as the time is wasted in producing candidate every time and computational cost is also more.

1.2 FP-growth

FP-growth algorithm is used for Frequent Itemset Mining. It is a depth-first search algorithm, which requires no candidate generation. FP-growth is faster than Apriori. In the mining process of FP-growth, there is not a single chance to re-truncate transactions. So, the transaction truncating approach is not suitable for FP-growth. In FP-growth during mining process, it is hard to obtain the exact number of support computations of i-itemsets. In FP-growth, the structure is an extended prefix-tree structure which is used for storing compressed and important information about frequent patterns.

FP-growth uses two data structures i.e. header table and FP-tree. Branch in FP-tree represents an itemset. Each node in FP-tree contains of a counter. In Header table, it stores items and their supports.

- Construction of the FP-Tree.
- Extracts frequent itemsets directly from the FP-Tree.

The FP-growth also referred as Private FP-growth (PFPgrowth) Algorithm. It consist of two modules preprocessing and mining. The preprocessing module is used to improve privacy as well as smart splitting method to transform the database. In the mining module a run time estimation method is used to estimate the support of item and dynamic reduction method to dynamically reduce the amount of noise. In between the formal privacy analysis which shows that the PFP-growth algorithm is ϵ -differentially private.

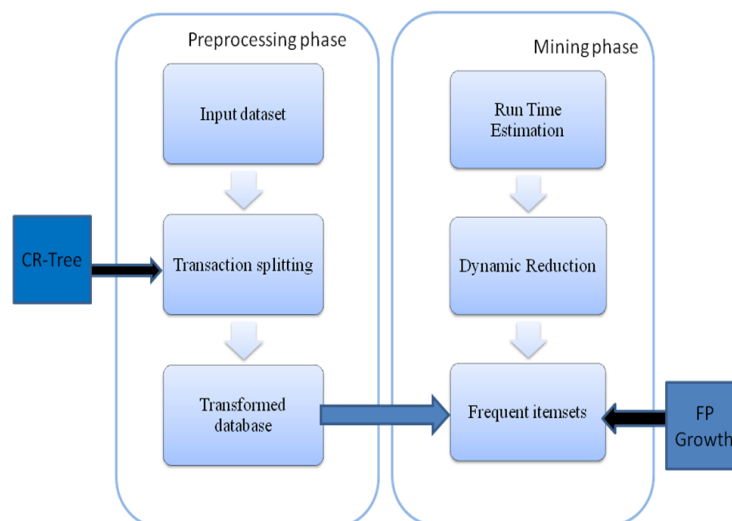


Fig.1. Existing system architecture with both phases [This diagram is pictorial representation as studied by authors]

1.3 Problem statement:-

- The existing system have many problem one of is to find with the high utility transactional itemsets.
- Existing methods takes more times for mining.
- Existing system gives many more combinations of output.
- Deleting the items until its length is under the limit, if a transaction has more items than the limit means its enforcing the limit by truncating the transactions.

IV. Analysis

Thus we have studied the ways from which we collect frequent itemset from long transaction of given database. But it doesn't consider the unfrequent but valuable itemset from given transaction. So, in Future we will work on utility itemset which is useful and profitable. For privacy preservation, in this both mining we uses differentially privacy method which gives more secure and private itemset. For private utility item mining will use UP-growth which gives us useful and profitable itemset as well as efficient time.

V. Conclusion And Future Work

In this paper, by survey so many methods for frequent item mining with privacy such as K-anonymity-diversity, Privbasis and also studies many different algorithms. Here, we examine that frequent itemset mining only consider the frequently occurring itemset and is challenged in many areas such as retail, marketing etc. It has been seen that in many real application domains that itemsets which share the most are not well enough the frequent itemset. So, in Future we will work on utility itemset which gives us useful and profitable itemsets as well.

REFERENCES

- [1] Li Zhou, Zhaohui Yang and Qing Yuan "Salient Region Detection via Integrating Diffusion-Based Compactness and Local Contrast," *IEEE Trans. Image processing*, vol. 24, no. 11, Nov. 2015
- [2] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention for rapid scene analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 11, pp. 1254–1259, Nov. 1998.
- [3] J. Harel, C. Koch, and P. Perona, "Graph-based visual saliency," in *Proc. Adv. Neural Inf. Process. Syst.*, 2006, pp. 545–552.
- [4] Y. Zhai and M. Shah, "Visual attention detection in video sequences using spatiotemporal cues," in *Proc. 14th Annu. ACM Int. Conf. Multimedia*, 2006, pp. 815–824.
- [5] X. Hou and L. Zhang, "Saliency detection: A spectral residual approach," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2007, pp. 1–8.
- [6] R. Achanta, S. Hemami, F. Estrada, and S. Süsstrunk, "Frequency-tuned salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 1597–1604.
- [7] S. Goferman, L. Zelnik-Manor, and A. Tal, "Context-aware saliency detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 2376–2383.
- [8] M.-M. Cheng, G.-X. Zhang, N. J. Mitra, X. Huang, and S.-M. Hu, "Global contrast based salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2011, pp. 409–416.
- [9] C. Lang, G. Liu, J. Yu, and S. Yan, "Saliency detection by multitask sparsity pursuit," *IEEE Trans. Image Process.*, vol. 21, no. 3, pp. 1327–1338, Mar. 2012
- [10] G.-K. Zhu, Q. Wang, and Y. Yuan, "Tag-Saliency: Combining bottom-up and top-down information for saliency detection," *Comput. Vis. Image Understand.*, vol. 118, pp. 40–49, Jan. 2014.
- [11] V. Gopalakrishnan, Y. Hu, and D. Rajan, "Salient region detection by modeling distributions of color and orientation," *IEEE Trans. Multimedia*, vol. 11, no. 5, pp. 892–905, Aug. 2009.
- [12] F. Perazzi, P. Krahenbuhl, Y. Pritch, and A. Hornung, "Saliency filters: Contrast based filtering for salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2012, pp. 733–740.
- [13] K. Shi, K. Wang, J. Lu, and L. Lin, "PISA: Pixelwise image saliency by aggregating complementary appearance contrast measures with spatial priors," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2013, pp. 2115–2122.
- [14] M.-M. Cheng, G.-X. Zhang, N. J. Mitra, X. Huang, and S.-M. Hu, "Global contrast based salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2011, pp. 409–416