# Improving Frequent Itemset Mining by Bifurcating Transactions

Aakash P Joshi, Prof. P. M. Kamde

PG Student, Department of Computer Engineering, Sinhgad College of Engineering, Pune, India

Associate Professor, Department of Computer Engineering Sinhgad College of Engineering, Pune, India.

**ABSTRACT**: There is need to protect the data and there has been huge rising curiosity in differential private data-mining algorithm. Differential privacy plans to give intends to amplify the exactness of questions from measurable databases while minimizing the odds of recognizing its records. Frequent itemset mining (FIM) is one of the most primary problems in data mining. In the existing system a differentially private FIM algorithm proposed which is based on the FP-growth algorithm mentioned as PFP-growth (Private- FP). The PFP-algorithm made up of two phases, pre-processing and mining phase. In the Pre-processing phase, smart splitting method is used to transform the database. In the mining phase, a run-time estimation method is used to measure the actual support of itemsets in the original database. But in smart splitting method, splitting transitions are sent one by one to the mining phase requires high memory and time to the system. Therefore in this paper the PFP-growth (Parallel-FP) is proposed which overcomes the memory and time problem. It uses the Thread instead of splitting and thread executes the operation in parallel.

**KEYWORDS**: Data-mining, information loss, transaction splitting and frequent itemset mining.

## I. INTRODUCTION

Frequent itemset mining plays an essential role in many data mining tasks that try to find interesting patterns from databases. The mining of data is one of the most popular problems of all these. The distinguishing proof of sets of things, items, side effects and qualities which regularly happen together in the given database, can be seen as a standout amongst the most essential undertakings in Data Mining.

The original motivation for searching frequent itemset came from the need to analyze so called supermarket transaction data that is to examine customer behavior. In existing system, private FP-growth (PFP-development) algorithm presented which comprised of a preprocessing phase and a mining phase. In the preprocessing phase, Database is changed to restrict the length of exchanges. The preprocessing stage is limited to user specified limits and should be performed just once for a given database. They argue to insist on this type of limit, lengthy transaction should be split i.e. from sub-transactions rather than truncate. For this a Smart Splitting technique is used to split the transactions.

In the Mining phase, the operation is performed on transformed dataset which is output of preprocessing phase. In this phase, Run-Time-Estimation and Dynamic Reduction technique is used. Run-Time-Estimation technique applied to find out information loss during mining phase. For maintain the privacy needs to add some amount of noise in the transactions. For finding the actual support of transactions (Original Database) and final support of transaction, Run-Time-Estimation is used. Dynamic Reduction technique is used to remove the noisy items in the transaction at the final stage i.e. in the mining phase after performing the run-Time-Estimation.

In Smart Splitting technique, after splitting the long transactions into sub-transactions it sends that transaction and their sub-transactions to mining phase one-by-one. So it takes high time to complete an operation and occupies memory to store all transactions for mining phase. To overcome this, we propose the Parallel FP-Growth algorithm. The PFP consists of five different steps, (Sharding, Parallel counting, Grouping Items, Parallel FP-Growth and Aggregation). It performs the operation by dividing the dataset using the Sharding method and assigns to the different-different threads and parallel Mining of the itemsets is done. In existing system we send the truncated translation one-by-one. It overcomes this by sending transactions in parallel. At the end of step grouping done, FP-tree algorithm is generated. In

aggregation step, gather all the items sets again to generate the frequent patterns.In this paper, further we will see: Section II talks about related work studied till now on topic. Section III includes current implementation details, introductory definitions.

## II. RELATED WORK

In this section discuss existing work done by the researchers for text mining process.In paper [1], to investigate the matter of planning a differentially private FIM algorithm presented. Private FP-growth (PFP-growth) algorithm was proposed that consisted of a preprocessing phase and a mining part. Within the preprocessing phase, to improve the utility-privacy tradeoff, conceived a wise rending methodology to remodel the information. In the mining part, a run-time estimation methodology planned to offset the knowledge loss incurred by group action splitting. Moreover, by using the downward closure property, dynamic reduction methodology applied to dynamically cut back the quantity of noises other to ensure privacy throughout the mining method. Formal privacy analysis and the results of intensive experiments on real datasets show that our PFP-growth algorithm is time-efficient and can achieve each sensible utility and privacy.

In paper [2], explained about parallel FP-Growth algorithm. The formula was predicated on a completely unique knowledge and computation distribution theme that eliminates communication among computers and makes it attainable for us to specify the formula with the Map-Reduce model. Experiments can be done on a vast dataset will give outstanding scalability of this formula. To form the formula suitable for mining internet knowledge, that are sometimes of long tail distribution, a tendency is to designed this formula to mine top-k patterns associated with every item, instead of wishing on a user specified worth for international smallest support threshold. Demonstration is done such that PFP is effective in mining associations and Webpage associations to support query recommendation or connected search.

In paper [3], it showed that a k-anonymity dataset permits solid attacks unpaid to lack of diversity within the sensitive attributes. Author introduced ℓ-diversity, a framework that gives stronger privacy guarantees. There square measure many avenues for future work. First is, to extend initial concepts for handling multiple sensitive attributes, and to develop ways for continuous sensitive attributes. Second, though privacy and utility square measures duals of alternative privacy, it has received far more attention than the utility of a broadcast table. As a result, the concept of utility wasn't well-understood.

In [4], author displayed two new algorithms, Apriori and Apriori-Tid, for finding all significant affiliation rules between things in a substantial database of exchanges. Author contrasted these algorithms with the beforehand known algorithms, the AIS and SETM algorithms. Author exhibited trial results, appearing that the proposed algorithms dependably outflank AIS what's more, SETM. The execution hole expanded with the issue measure, and ran from an element of three for little issues to more than a request of greatness for huge issues.

In paper [5] Sequential data is progressively utilized in many application. Publishing sequential data is of basic importance to the advancement of these applications. However, as appeared by the re-identification attacks on the AOL and Netix datasets, discharging sequential data might posture extensive dangers to individual privacy. Differential privacy is one of the main models that can be utilized to give such ensures. Our approach is to design tree structure and set of novel techniques depends on the Markov assumption for minimizing the added noise . The distributed n - grams are valuable for some reasons. Besides, we add to an answer for producing an synthetic database, which empowers a more wider spectrum of data analysis tasks. Broad tests on real-life datasets show that our methodology significantly outperforms the state-of-the-art techniques.

A. System Overview

The following figure 1 shows the architectural view of the proposed system. The description of the system is as follows:

In the proposed system initially input is a Dataset file. Input is given to the Parallel FP-Growth (PFP) algorithm .For the Database DB,PFP utilizes three Map-Reduce stages to parallelize PF-Growth.
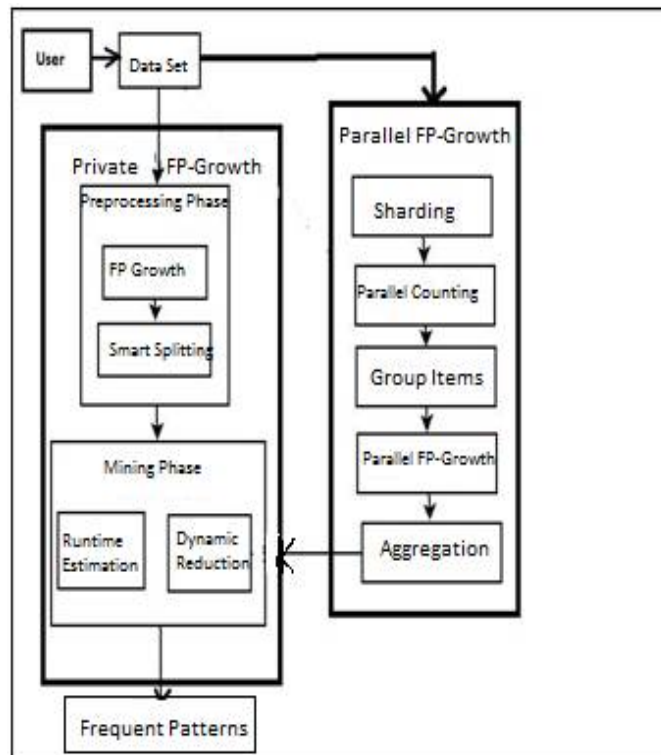


Figure 1: System Architecture

B.   Algorithm

Proposed System Algorithm
Algorithm 1: The Parallel Counting Algorithm

**Input:** Dataset.
**Output**: Divided dataset.

**Process**:
Step 1:  Input the Dataset.
Step 2:  Sharding the dataset in number of shard.
Step 3:  Assign each shard to each different Thread.
Step 4: Perform parallel counting.
Step 5: Grouping Items.

Algorithm 2: The Parallel FP-Growth Algorithm

**Input:** Divided dataset.

**Output**: Frequent Itemset in each thread.

**Process**:
Step 1: Generate the FP-tree
Step 2: Generate the header table

Step 3: Create a Mapper.
Step 4: Generate Hash Table *H* from G-List.
Step 5: Reducer

Algorithm 3: The Aggregating Algorithm

**Input:** Frequent Itemset in each thread.
**Output**: Frequents Patterns.

**Process**:
Step 1: Getting all threads.
Step 2: Create a Mapper
Step 3: Generate Hash Table *H* from G-List
Step 4: Reducer

## C. PROPOSED ALGORITHM

System S is represented as S = {I, F, N, U, L, T, K, P, R, D, O }

**Input:**
I = {D, N} Where D = Sparse datasets:  Retail Dataset a

   N = {n1, n2, n3, n4 , n5}

Where, N is the set of Noise and n1, n2, n3, n4 , n5 are the noise sample.

**Process:**

- F = FP Growth algorithm
- Noise
  N = n1
  Where N is a noise and n1 is a noise sample.

- Undirected Weighted Graph U = { u1, u2, u3,....,un }

Where U is represent as a graph and u1, u2, u3, .....,un
Number of vertices from the graph
- L= Louvain Method
- T= CR-Tree

$$\Delta Q = \left[ \frac{\sum_{in} + k_{i,in}}{2m} - \left( \frac{\sum_{tot} + k_i}{2m} \right)^2 \right] -$$

$$\left[ \frac{\sum_{in}}{2m} - \left( \frac{\sum_{tot}}{2m} \right)^2 - \left( \frac{k_i}{2m} \right) \right]$$

Where,
$\Delta Q$ = used to find the Modularity of node in CR-Tree.

$\sum_{in}$ = total of the considerable number of weights of the links inside the community i is moving into.

$\sum_{tot}$ = som of number of weights of the links to nodes in the community.

**k_i** = degree of i.

**k_{i, in}** = total of the weights of the links in the middle of i and different nodes in the community.

and

**m** = whole of the weights of all links in the networks.

- Splitting K = { k1, k2, k3, ....kn }

Where K be a set of Splitting where, k1, k2, k3,......,kn are the number of split tree.

- P= Parallel FP-Growth
- R= Run-time Estimation

$$\omega_a = \mathrm{avg\_supp}(\varpi, i) = \int_{\omega'=\varpi-5}^{\omega'=\varpi+5} \mathrm{Pr}(\omega' \mid \varpi) avg(\omega', i)$$

For calculating Average Support of ItemSet ,we have

$$\omega_m = \mathrm{max\_supp}(\varpi, i) = \int_{\omega'=\varpi-5}^{\omega'=\varpi+5} \mathrm{Pr}(\omega' \mid \varpi) \max(\omega', i)$$

For calculating Maximum Support of Item Set, we have,

- D= Dynamic Reduction

**Output:**

O = Reduce noise and Privacy with frequent Patterns

#### D. Experimental Setup

The system is built using Java framework (version jdk 8) on Windows platform. The Netbeans (version 8.1) is used as a development tool. The system doesn't require any specific hardware to run; any standard machine is capable of running the application.

### IV. RESULT AND DISCUSSION

a. DataSet

Sparse datasets: Retail Dataset. This dataset contain numeric value for items. Contain 200 transaction of different length.

b. Results

In this paper we find the frequent item-sets. In existing system, we send the truncated transactions one-by-one. In proposed, mining is done parallel. It's sparse the time as well as memory.

Figure 2:  Reading Dataset



Figure 3: Splitting done in existing System

Figure 4: Runtime estimation method



Figure 5: Dataset Division in Proposed Parrallel FP growth Algorithm.
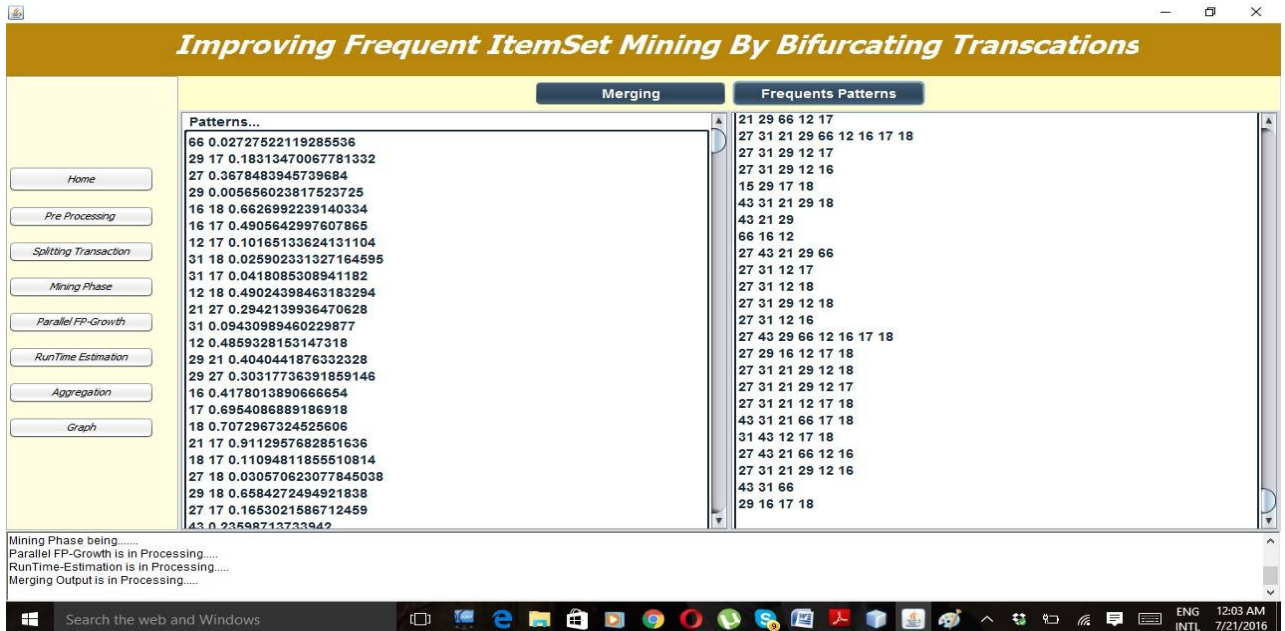
Figure 6: Aggregation and final frequent patterns.

For different value of threshold, estimated result show's that, the memory and time require for the proposed system is less than existing system.
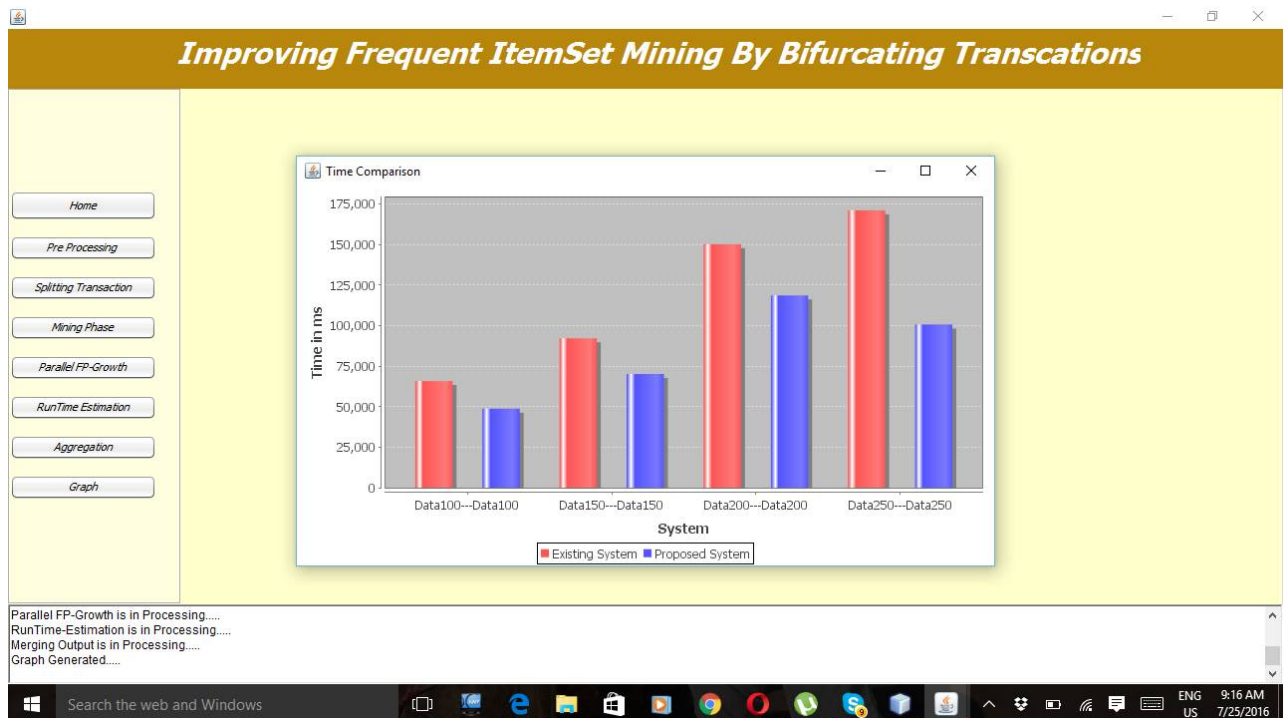


Figure 7: Time Comparison graph.

Conclusion and Future scope

In this paper, we explore the issue of planning a differentially private FIM algorithm. We propose our parallel FP-growth (PFP-growth) algorithm, which consists of a five different steps. PFP-algorithm is time-productive and require less memory.
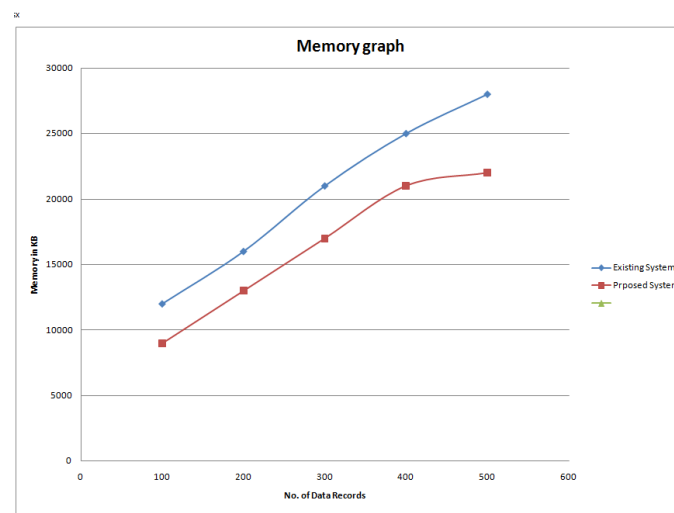


Figure 8: Memory Comparison graph.

### REFERENCES

[1]Sen Su, Shengzhi Xu, Xiang Cheng, Zhengyi Li, and Fangchun Yang, "Differentially Private Frequent Itemset Mining via Transaction Splitting", IEEE Transactions on Knowledge and Data Engineering,DOI10.1109/ TKDE.2015.2399310 ,2015.
[2] Haoyuan Li, Yi Wang, Dong Zhang, Ming Zhang, Edward Chang," PFP: Parallel FP-Growth for Query Recommendation ", Google Beijing Research, Beijing, 100084, China.
[3] Ashwin Machanavajjhala ,Johannes Gehrke, Daniel Kifer,Muthuramakrishnan Venkitasubramaniam, "ℓ-Diversity: Privacy Beyond k-Anonymity", Department of Computer Science, Cornell University.
[4] Rakesh Agrawal, Ramakrishnan Srikant, "Fast Algorithms for Mining Association Rules", IBM Alma den Research Center 650 Harry Road, San Jose, CA 95120
[5] R. Chen, G. Acs, and C. Castelluccia, "Differentially Private Sequential Data Publication via Variable-Length n-Grams," Proceedings of the 2012 ACM conference on Computer and communications security. ACM , 2012.
[6] L. Bonomi and L. Xiong, "A Two-Phase Algorithm for Mining Sequential Patterns with Differential privacy," CIKM, 2013.
[7] E. Shen and T. Yu, "Mining Frequent Graph Patterns withDifferential Privacy," in KDD, 2013.
[8]M. Hay, V. Rastogi, G. Miklau, and D. Suciu. "Boosting the Accuracy of Differentially Private Histograms through Consistency". VLDB, 3(1):1021-1032, 2010.
[9] X. Xiao, G.Wang, and J. Gehrke. "Differential privacy via Wavelet Transforms", IEEE Transaction on Knowledge. Data Eng., 23(8):1200-1214, 2011.
[10]Chen Zeng and Jin-yi Cai, "On Differentially Private Frequent Itemset Mining", Proceedings of VLDB Endowment, Vol 6, NO. 1, 26th august- 30th 2013.
11]Raghav Bhaskar, Srivatsan Laxman, "Discovering Frequent Patterns in Sensitive Data" ACM KDD, 2010, Washington DC.
[12]Ninghui Li, W. Qardaji, Dong Su, Jianneng Cao,"PrivBasis: Frequent Itemset Mining with Differential Privacy " Proceedings of the VLDB Endowment, Vol 5,No.11, August 27th-31st 2012 Istanbul Turkey.

## BIOGRAPHY

**Aakash Joshi** is currently pursuing M.E(Computer Networks) from Department of Computer Engineering, Sinhgad College of Engineering, Pune,Savitribai Phule Pune University, Pune, Maharashtra, India-411041. He received his B.E (Computer Engineering) Degree in 2013 from Trinity College of Engineering,Pune, Savitribai Phule Pune University, Pune, Maharashtra, India -411048. His area of interest is Data Mining and Information Security.

**Prof. Pravin M. Kamde** is a Associate Professor in the Department of Computer Engineering at Sinhgad College of Engineering,Pune,India. He has received his degree of B.E.(Computer Science and Engineering) from SGGS College of Engineering and Technology, Nanded,Marathwada University of Aurangabad in 1993,M.E.(Computer Science and Engineering) in 2004 from Walchand College of Engineering, Sangli, Shivaji University. He has 13 Internal Journal Publications, 7 International Conferences and 14 National Conferences publication. his research interests include Content-Based Image and video Retrievaal, Web Multimedia Mining and Image Processing.