# Frequent Itemset mining based on Differential Privacy using RElim(DP-RElim)

## Bhagyashree R. Vhatkar[1], Prof. Dr. Mrs.S.A.Itkar[2]

[1]Department of Computer Engineering P.E.S. Modern College of Engineering,Pune
[2]Department of Computer Engineering P.E.S. Modern College of Engineering,Pune,suhasini_

-------------------------------------------------------------------------***-------------------------------------------------------------------------

**Abstract -** In recent years, individuals are interested in designing differentially private data mining algorithms. Many researchers are working on design of data mining algorithms which gives differential privacy. In this paper, to explore the likelihood of planning a differentially private FIM , cannot just accomplish high information utility and a high level of protection, additionally offers high time effectiveness. To this end, the differentially private FIM based on the FP-growth algorithm, which is speak about to as PFP-growth. The Private RElim algorithmic program consists of a pre-processing part and a mining part. within the preprocessing part, to enhance the utility and privacy exchange, a completely unique good smart splitting technique is expected to rework the database.A frequent itemset miningwith differential privacy is important which will follow twophase process of pre-processing and mining. Through formal private investigation, demonstrate that our Private DP-RElim is "ε- differentially private. Broad analyses on genuine datasets show that our DP-RElim algorithm considerably outflanks the best in class systems.The computational experiments on real world and synthetic databases exhibit the fact that in comparison to the performance of previous algorithms, our algorithms are faster and also maintain high degree of privacy, high utility and high time efficiency simultaneously.

*Key Words***:** Frequent itemset mining, Differentially private, Pre-processing, Mining, Private RElim, Transaction splitting.

## 1.INTRODUCTION

In the database, where every exchange contains an arrangement of things, FIM tries to discover item-sets that happen in exchanges more much of the time than a given limit. An assortment of algorithms have been proposed for mining incessant itemsets. The Apriori what's more, FP-algorithm are the twomost essential ones. Specifically, Apriori is a broadness first pursuit, competitor set era andtest lgorithm. It needs one database examines if the maximal length of incessant itemsets is one . Conversely, FP-growth is a profundity first hunt algorithm which requires no applicant era. In FP-growth just performs two database checks, which makes Frequent Pattern a request of greatness speedier than Apriori. The engaging components of FP-growth inspire us to outline a differentially private FIM

algorithm in light of the FP algorithm. In this paper, the differentially private FIM ought not just accomplish high information utility and a high level of security, additionally offer high time productivity. Although a few differentially private FIM algorithms have been proposed, they don't know about any current reviews that can fulfill every one of these necessities all the while. The subsequent requests fundamentally bring new difficulties.

In past work shows an Apriori-base ddifferentially private FIM algorithm. It implements the breaking point by truncating. In specific, in every database sweep, to safeguard more recurrence data, it favourable position to found regular itemsets to re-truncate exchanges. Nonetheless, FP-growth just performs two database checks. There is no chance to retruncate exchanges amid the mining procedure. Subsequently, the exchange truncating methodology is not reasonable for FP-growth . Furthermore, to maintain a strategic distance from security break, the add commotion to the support of itemsets. FP-growth is a profundity first inquiry algorithm not like Apriori. It is difficult to get the correct number of bolster algorithms of i-itemsets amid the mining procedure. An innocent way to deal with figure the boisterous supportof i-itemset isto utilize the quantity of all conceivable i-itemsets. In any case, it will certainly create invalid outcomes. Apriori-based is significantly upgrade by transaction splitting techniques:

-The return to the exchange off amongst utility and security in outlining a differentially private FIM . The exhibit that the exchange off can be expanded by our novel transactionsplitting techniques exchange part procedures. Such procedures are appropriate for FP-growth, as well as can be used to plan other differentially private FIM .

-To create a period effective differentially private FIM algorithm in light of the FP-growth algorithm which is alluded to as PFP-growth. Specifically, by utilizing the descending closureproperty, a dynamic lessening technique is proposed to progressively lessen the measure of commotion added to ensure security amid the mining procedure.

- Through formal privacy investigation, the demonstrate that our PFPgrowth algorithm is "- differentially private.

The subsequent sections of the paper are organized as follows.Section II gives brief review of the differential private techniques for databases.Section III provides with a detailed

information about the system architecture. Section IV gives the experimental analysis and performance parameter considerations. In Section V,the results for implemented system is analyzed.Finally.the paper is concluded in Section VI.

## 2. REVIEW OF LITERATURE

Sen Su, Shengzhi Xu[1].In this paper the components of FP-growth propel us to plan a differentially private FIM algorithm in view of ther FP-growth algorithm. We contend that a reasonable differentially private FIM algorithm ought not just accomplish high information utility and high degree of security, additionally offer high time effectiveness. FP-growth just performs two database check. There is no open door to re-truncate exchange amid mining process. Private FPgrowth(PFP-growth) algorithms, which comprise of preprocessing stage and mining stage. In preprocessing stage we change the database to restrain the length of exchanges. The preprocessing stage is superfluous to userspecified edges furthermore, should be performed once for a given database. That is, if an exchange has a greater number of things than the point of confinement, we separate it into various subsets and ensure every subset is under the breaking point. We devise a smart splitting strategy to change the database. Specifically, to guarantee applying - differentially private algorithm on the changed databasestill fulfills ε- differential protection for the unique database, we propose a weighted splitting operation. Additionally, to more recurrence data insubsets, we propose a graph based way to deal with uncover the relationship of things inside exchanges and use such correlationto direct the splitting procedure.

Zeng C[2]In this paper,we concentrate on security issues that emerge with regards to finding continuous itemsets in value-based information. It can investigate the likelihood of growing differentially private incessant itemset mining algorithms. We will likely ensure differential protection without destroying the utility of the algorithm. A nearer examination of this negative outcome uncovers that it depends on the likelihood of long exchanges. This raises the likelihood of enhancing the utility-protection exchange off by restricting exchanges cardinality. Obviously, one can't as a rule force such a farthest point; so all things being equal, we investigate upholding the farthest point by truncating exchanges [2]. That is, if an exchange has more than a predefined number of things, we erase things until the exchange is under the breaking point. Obviously, this cancellation must be done in a respectfully private manner; maybe similarly essential, while it diminishes the mistake due to the commotion required to authorize security. Thought of restricting the maximal cardinality oftransactions is basic we truncate an exchange whose cardinality damages that requirement by just keeping a subset of that exchange. Obviously, that truncating approach brings about certain data misfortune. However,if the cardinality of exchanges in a

dataset takes after a dispersion in which most are short and a couple are long, then these few long exchanges, while having little effect on which itemsets are visit, have a noteworthy impact on the sensitivity.

Ninghui Li[3]. In this paper rxamine a novel approach that maintains a strategic distance from the determination of top k item-sets from a large competitor set. All the more uncommonly, we present the thought of premise sets. A θ-premise set(B) = B1;B2;Bw; where everyBi is an arrangement of things, has the property that anyitemset with recurrence higher than θis a subsetof some premise Bi. A decent premise setis one where w is little and the lengths of all Bi are additionally little. Given a decent premise set B, one can remake the frequencies ofall subsets of Bi with great precision. One can then choose the most incessant itemsets from these. We additionally acquaint systems with build great premise sets while fulfilling differential protection. It meets the test of high dimensionality by anticipating the information informational collection onto a little number of chose measurements that one thinks about. Actually, PrivBasis regularly utilizes a few arrangements of measurements for such projections, to stay away from any one set containing an excessive number of measurements. Every premise in B relates to one such arrangement of measurements for projection. Our strategies empower one to choose which sets of measurements are most useful with the end goal of discovering thek most incessant itemsets. A key idea presented in this approach is the idea of Truncated Frequencies (TF). The TF strategy tries to address the running time challenge by pruning the pursuit space, however it doesn't address the exactness challenge.

J.han,J.pei[4],Mining regular examples in exchange databases, time-arrangement databases, and numerous different sorts of databases has been contemplated prevalently in information mining research. The majority of the past reviews receive an Apriori-like hopeful set era and-test approach. Notwithstanding, hopeful set era is still exorbitant, particularly when there exist an extensive number of examples and additionally long examples. In this review, to propose a novel regular example tree (FP-tree) structure, which is an expanded prefix-tree structure for securing compacted, fundamental information about normal cases and build up an productive FP-tree based mining strategy, FP-growth, for mining the entire arrangement of regular examples by example section growth. Productivity of mining is accomplished with three procedures: (1) an extensive database is packed into a dense, littler information structure, FP-tree which dodges exorbitant, rehashed database examines, (2) our FP-tree-based mining receives an example piece growth strategy to stay away from the exorbitant era of a vast number of hopeful sets, and (3) an apportioning based, divideand- conquer technique is utilized to break down the mining undertaking into an arrangement of littler errands for mining bound examples in contingent databases, which significantly diminishes the hunt space

Apriori algorithm and furthermore speedier than some as of late revealed new regular example mining techniques.

Vaidya and C.Clifton[5],This paper addresses the issue of affiliation govern mining where exchanges ar circulated crosswise over sources. Every site holds a few characteristics of every exchange, what's more, the locales wish to work together to distinguish all inclusive legitimate affiliation rules. Be that as it may, the locales must not uncover person exchange information. We display a two-party calculation for effectively finding continuous itemsets with least support levels, without either site uncovering singular exchange values. To introduce a system for mining affiliation rules from exchanges comprising of categorical things where the information has been randomized to protect security of individual exchanges. While it is practical to recuperate affiliation principles and protect protection utilizing a clear uniform randomization, the found principles can lamentably be misused to and security ruptures examine the nature of security breaks and propose a class of randomization administrators that are a great deal more compelling than uniform randomization in constraining the ruptures. the determine formulae for a fair-minded bolster estimator and its fluctuation, which permit us recoup itemset bolsters fro randomized datasets, and demonstrate to join these formulae into mining calculations. At last, to show test comes about that approves the calculation by applying it on genuine datasets.By vertically parceled, the imply that every site contains a few components of an exchange. Utilizing the customary market crate case, one site may contain basic supply buys, while another has dress buys. Utilizing a key, for example, credit card number what's more, date, it can join these to distinguish connections between buys of dress and staple goods. Nonetheless, this reveals the singular buys at every site, perhaps abusing shopper protection assentions. There are more sensible illustrations. In the sub-gathering fabricating process, distinctive makers give parts of the completed item. Autos fuse a few subcomponents; tires, electrical gear, etc.made by autonomous makers.

Bhaskar R[6],In this paper exhibit two productive calculations for finding the K most incessant examples in an informational collection of delicate records. Our calculations fulfill differential security, an as of late presented definition that gives important protection ensures within the sight of discretionary external information. Differentially private calculations require a level of vulnerability in their out-put to protect security. Our calculations handle this by returning uproarious arrangements of examples that are near the genuine Run down of K most continuous examples in the information. We characterize another idea of utility that evaluates the yield exactness of private best K design mining algorithms.[6]

L.Bonomi[16] Visit successive example mining is a focal undertaking in many fields, for example, science and back. In any case, arrival of these examples is raising expanding worries on individual security. In this paper, concentrate the successive example mining issue under the differential security system which gives formal and provable certifications of security. Because of the way of the differential protection component which bothers the recurrence comes about with commotion, and the high dimensionality of the example space, this mining issue is especially testing. In this work, the propose a novel two-stage calculation for mining both prefixes what's more, substring examples. In the principal stage, our approach takes favorable position of the measurable properties of the information to build a model-based prefix tree which is utilized to mine prefixes and a competitor set of substring examples. The recurrence of the substring examples is further refined in the progressive stage where the utilize a novel change of the first information to decrease the annoyance commotion.

Christian Borgelt[17] In this paper a recursive elimination scheme: in a preprocessing step delete all items from the transactions that are not frequent individually,i.e., do not appear in a user-specified minimum number of transactions. Then select all transactions that contain the least frequent item, delete this item from them, and recurse to process the obtained reduced database, remembering that the item sets found in the recursion share the item as a prefix.On return, remove the processed item also from the database of all transactions and start over, i.e., process the second frequent item etc. In these processing steps the prefix tree, which is enhanced by links between the branches, is exploited to quickly find the transactions containing a given item and also to remove this item from the transactions after it has been processed.It processes the transactions directly, organizing them merely into singly linked lists. The main advantage of such an approach is that the needed data structures are very simple and that no re-representation of the transactions is necessary, which saves memory in the recursion. In addition, processing the transactions is almost trivial and can be coded in a single recursive function with relatively few lines of code. Surprisingly enough, the price one has to pay for this simplicity is relatively small: my implementation of this recursive elimination scheme yields competitive execution times.

Christian Borgelt [18] in this paper the RElim (Recursive Elimination) algorithm can be seen as a precursor of the SaM algorithm. It also employs a basically horizontal transaction representation, but separates the transactions (or transaction suffixes)according to their leading item, thus introducing a vertical representation aspect.In addition, the transactions are organized as lists (at least in my implementation), even though,in principle, using arrays would also be possible. These lists are sorted descendingly w.r.t. the frequency of their associated items in the transaction database: the first list is associated with the most frequent item, the last list with the least frequent item.

## 3. SYSTEM ARCHITECTURE / SYSTEM OVERVIEW

The DP-RElim algorithm consists of a preprocessing phase and a mining phase. In the preprocessing phase, a novel smart splitting method is proposed to transform the database. In the mining phase, a run-time estimation method to estimate the actual support of itemsets in the original database, we put forward a dynamic reduction method to dynamically reduce the amount of noise.
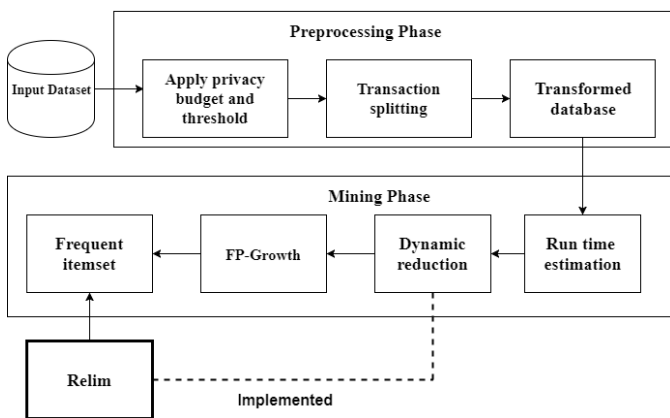


**Fig -1:** Architecture diagram

## 3.1 Processing steps:

### 1) Input Data Collection

We collect two dataset from http://fimi.ua.ac.be/data/ first traffic accident data and second contains the retail marketbasket data.

### 2) ε-Differential privacy

By add a rigorously chosen quantity of noise, differential privacy assures that the output of a estimation is insensitive to changes in any people record, and so limiting privacy leaks through the results.

### 3) Transaction Splitting

To limit the cardinality of transactions by transaction splitting, we can keep more frequency information. That is, long transactions are divided into multiple sub transactions whose cardinality is below a specified number of items.

### 4) Transformed database

The database to restrict the length of transactions. To uphold such a breaking point, long exchanges ought to be part instead of truncated.

### 5) Run time estimation
A run-time estimation method is proposed to balanced

the information loss obtain by transaction splitting.

### 6) Dynamic reduction

To dynamically decrease the amount of noise added to ensure privacy amid of mining process.

### 7) FP-Growth

FP-growth is a depth first search algorithm, which requires no candidate generation.

### 8) Recursive Elimination(Implemented)

Recursive Elimination algorithm relies on a step by step elimination of items from the transaction database together with a recursive processing of transaction subsets This algorithm works without complicated data structures and allows us to find frequent item-set easily.

## 4. EXPERIMENTAL SYSTEM

### 4.1. Hardware and software Requirement

Experiments are conducted on Processor: Intel Duo Core2 E8400 CPU(2.0 GHz)and 4 GB RAM, HDD: 1 TB System Type: 64 Bit.The front end used is JAVA jdk 1.8 with Netbeans IDE 8.0.2 and the system is implemented on Windows 7 operating system. In the experiments, we use two publicly available real datasets. Retail dataset which contain market basket data and Accident dataset which contain traffic accident data.

### 4.2. Performance parameter

To calculate the performance of algorithm, we utilize the widely used standard metrics.

1) F-score: It measures the utility of generated frequent itemset.

$$F - score = 2 * \frac{precision * recall}{precision + recall}$$

where,

$$Precision = \frac{|U_p \cap U_c|}{U_p}$$

$$recall = \frac{|U_p \cap U_c|}{U_c}$$

$U_p$ is frequent itemsets generated by private algorithm.

$U_c$ is the actual frequent itemset.

## 5. RESULTS

Table 1 F-Score in percentage with different threshold value on Retail dataset.

| Threshold | F-Score in % | |
| | Base paper | Implemented paper |
|---|---|---|
| 0.54 | 0.8 | 0.84 |
| 0.58 | 0.84 | 0.90 |
| 0.62 | 0.88 | 0.95 |

**Table 1-** F-Score in percentage with different threshold values on Retail dataset.

Graph gives value of F-score in percentage using DP-RElim. The x-axis consist of threshold value while Y-axis consist of F-score.
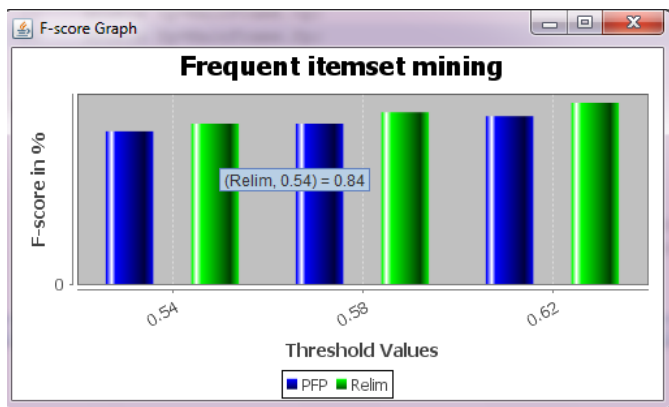


**Fig 2**. F-Score in % on Retail dataset.

Table 2:Running time evaluation(Time in milliseconds)for top-k frequent itemset on Retail dataset.

| Top-k Frequent itemset | Base paper | Implemented paper |
|---|---|---|
| 10 | 10000 | 6000 |
| 20 | 17000 | 11200 |
| 30 | 30000 | 24000 |

**Table 2.** Running time evaluation (Time in milliseconds) for top-k frequent itemset on Retail dataset.

Following graph running time evaluation.The x-axis consist of Top-k frequent itemset value while Y-axis consist of time in milliseconds.
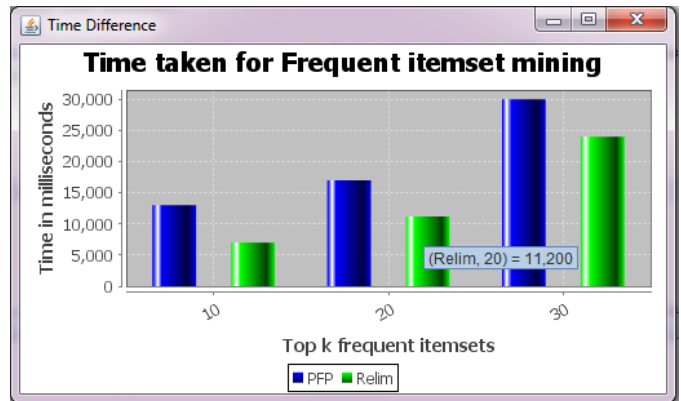


**Fig 3.** Running time evaluation on Retail dataset

## VI. CONCLUSION

In this paper we examine the problem of design a private DP-RElim with differential privacy ,which consist of preprocessing phase and mining phase. In first phase to better enhance utility exchange off, utilizing keen part strategy. In mining stage, a run time estimation technique is proposed to counterbalance the data misfortune brought about by exchange part. By using dynamic reduction method to dynamically decrease the amount of noise added to guarantee privacy during the mining process. The DP-RElim algorithm is time efficient and can achieve both utility and good privacy. The dynamic reduction and run-time estimation methods are used in phase to enhance the quality of the results. Recursive depends on a stage by step end of things from the exchange database together with a recursive preparing of exchange subsets. This calculation works without entangled information structures furthermore, permits us to discover visit itemset effectively.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Sen Su, Shengzhi Xu, Xiang Cheng, Zhengyi Li," Differentially Private Frequent Itemset Mining via Transaction Splitting" IEEE Transaction on knowledge and Data Engineering, vol. 27, No. 7, July 2015.

[2] C. Zeng, J. F. Naughton, and J.-Y. Cai, On differentially private frequent itemset mining," International Conference on Very Large Data Bases, Vol. 6, August 2012.

[3] N. Li, W. Qardaji, D. Su, and J. Cao, Privbasis: Frequent itemset mining with differential privacy", International Conference on Very Large Data Bases, Vol. 5, August 2012.

[4] J. Han, J. Pei, and Y. Yin, Mining frequent patterns without candidate generation, in SIGMOD, 2000.

[5] J. Vaidya and C. Clifton, Privacy preserving association rule mining in vertically partitioned data, in KDD,2002.

[6] Cynthia Dwork. "Differential Privacy" ICALP, Springer, 2006.

[7] R. Agrawal and R. Srikant, Fast algorithms for mining association rules, in Proc. 20th Int. Conf. Very Large Data Bases, 1994, pp. 487499.

[8] J. Han, J. Pei, and Y. Yin, Mining frequent patterns without candidate generation, in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2000.

[9] M. Kantarcioglu and C. Clifton, Privacy-preserving distributed mining of association rules on horizontally partitioned data, IEEE Transaction Knowledge Data Eng., vol. 16, no. 9, pp. 10261037, Sep. 2004.

[10] R. Chen, G. Acs, and C. Castelluccia, Differentially private sequential data publication via variable-length n grams, in CCS, 2012.

[11] Nandhini, Madhubala, Valampuri," Privacy-Preserving Private Frequent Itemset Mining via Smart Splitting", International Journal of Innovative Re-search in Computer and Communication Engineering, Vol. 3, October 2015.

[12] Anusuya M,Sudharani K,Ganthimathi M,Sumathi G," Frequent Itemset Mining Using PFP-Growth via Transaction Splitting", International Journal of Innovative Research in Computer and Communication Engineering, Vol. 4, Issue 2, February 2016.

[13] R. Bhaskar, S. Laxman, A. Smith, and A. Thakurta, Discovering frequent patterns in sensitive data, in Proc. 16th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2010.

[14] C. Dwork, F. McSherry, K. Nissim, and A. Smith, Calibrating noise to sensitivity in private data analysis, in Proc. 3rd Conf. Theory Cryptography, 2006.

[15] X. Zhang, X. Meng, and R. Chen, Differentially private set valued data release against incremental updates, in Proc. 18th Int. Conf. Database Syst. Adv. Appl., 2013.

[16] L. Bonomi and L. Xiong, A two-phase algorithm for mining sequential patterns with differential privacy, in Proc. 22nd ACMConf. Inf. Knowl. Manage., 2013.

[17] Christian Borgelt,"Simple Algorithms for Frequent Item Set Mining" Advances in Machine Learning II,Springer,2010.

[18] Christian Borgelt,"Keeping Things Simple: Finding Frequent Item Sets by Recursive Elimination"