# Frequent Item-set mining with differential privacy in Large Scale data

## Miss.Nikita DevramIlhe[1] Prof.ManojWakchaure[2]

*P.G.Student,DepartmentofComputerEngineeringAssistantProfessor,DepartmentofComputerEngineering,
AmrutvahiniCollegeofEngineering,Sangamner,MH,India[1]AmrutvahiniCollegeofEngineering,Sangamner,MH,India[2]*

***Abstract:*** *In recent years, people have an interest in coming up with differentially non-public data processing algorithms. several researchers square measure functioning on style of information mining algorithms which supplies differential privacy. during this paper, to explore the chance of coming up with a differentially non-public FIM , cannot simply accomplish high info utility and a high level of protection, in addition offers time effectiveness. to the present finish, the differential non-public FIM based on the FP-growth algorithmic program, that is talk about to as PFP-growth. The Privacy based mostly algorithm consists of a pre-processing half and a mining half. among the preprocessing half, to boost the utility and privacy exchange, a very distinctive smart good ripping technique is anticipated to transform the information. Privacy is the most essential in today's world for online as well as offline data. Frequent Itemsets Mining (FIM) it is a typical data processing task and has gained abundant attention. Due to the consideration of individual privacy, various studies have been focusing on privacy-preserving FIM problems. Differential privacy has emerged as a promising theme for shielding individual privacy in data processing against adversaries with impulsive information. In this  work we propose an efficient, privacy preservation based frequent itemsets mining (FRM) algorithm on large as well as high dimensional data called Frequent Itemset Mining Privacy Preservation (FIM_PP). In light of the thoughts of examining and exchange truncation utilizing length limitations, our calculation lessens the calculation force, decreases mining affectability, and subsequently enhances information utility given a settled protection spending plan. Partial experimental analysis show the proposed system evaluation show the how proposed system provides best results than existing systems.*

***Keywords:*** *Wearable sensors, healthcare, bigdata, cloud computing, authentication, security.*

## I. Introduction

In the info, wherever each exchange contains a meeting of things, FIM tries to get item-sets that happen in exchanges a lot of abundant of the time than a given limit. An assortment of algorithms is planned formining incessant itemsets. The Apriori what's a lot of, FPalgorithm area unit the twomost essential ones. Specifically, Apriori may be a breadth initial pursuit, competition set era andtest algorithm. It desires one info examines if themaximal length of incessant itemsets is one . Conversely, FPgrowth may be a profundity initial hunt algorithmic program which needs no human era. In FP-growth simply performs 2 info checks, that makes Frequent Pattern letter of invitation of greatnessspeedier than Apriori. The participating parts of FPgrowth inspire North American country to stipulate a differentially non-public FIM algorithmic program in lightweight of the FP algorithmic program. during this paper, the differentially non-public FIM ought not simply accomplish highinformation utility and a high level of security, to boot supply time productivity. though some differentially non-public FIM algorithms are planned, they don't understand any current reviews that may fulfill all of thosenecessities all the whereas. the following requests essentially bring new difficulties.

In past work shows AN Apriori-base differentially non-public FIM algorithmic program. It implements the brink by truncating. In specific, in each info sweep, to safeguard a lot of return knowledge, it favorable position to found regular itemsets to re-truncate exchanges. all the same, FP-growth simply performs to info checks. There's no likelihood to re-truncate exchanges amid the mining procedure.Subsequently, the exchange truncating methodology isn't affordable for FP-growth . moreover, to keep up a strategic distance from security break, the add commotion to the support of itemsets. FP-growth may be a profundity initial inquiry algorithmic program not like Apriori. it's tough to urge the proper variety of bolster algorithms of i-itemsets amid themining procedure. AN innocent thanks to handle figure the boisterous supportof k-itemset isto utilize the amount of all conceivable i-itemsets. In any case, it'll definitely produce invalid outcomes..

## II. Literature Survey

Pan, Zhaopeng et.al [1] FP-growth algorithm produces all the frequent item sets without producing a large number of candidate items. However, when the item set is too large, the branch of the spanning tree will be

long, occupying the space too large, and the mining efficiency will be reduced. In this paper, the method of using dynamic insert node FP - tree structure, and all the back pointer, to generate a new type of FP - tree. This article also proposes Max-IFP maximum frequent patterns mining algorithm, using the new generation of FP - tree dug up all the maximum frequent item sets. The experimental results show that the new FP-tree occupies a smaller space, and the algorithm proposed in this paper is shorter and more effective than other algorithms when mining the maximum frequent item sets.

Fahrudin, TresnaMaulana et.al [2] No proof of malady (NED) is carcinoma patient condition standing that it indicates that they will life, no realize the cancer by tested, and with none symptoms of cancer in period of times, after they received primary treatment. Patient condition status which it indicates that they NED is a critical status, because it involves the treatment type and patient cancer condition factors. Theexamines about breast cancer problem in data mining technical side, especially to discover the patterns of NED-breast cancer patient using cancer registry data from Oncology Hospital. Its patterns are discovered through the relationship of among features begin from 1dimensional, 2-dimensional, 3-dimensional, and n-dimensional. They applied association rules mining using Apriori and FP-Growth algorithm, which both have the advantage and drawback. Apriori algorithm involves all generation of candidate item sets and multiple database scans, but it makes high consuming iteration. While FP-Growth algorithm extracts the frequent item sets directly from FP-Tree, it make the advantage of FP-Growth that is faster process needs only scan the database once. This paper experiment shown that the association result of Apriori and FP-Growth is almost similar, 10-highest confidence value represented 100% confidence of association rule on breast cancer dataset with support value up to 50%.

Djenouri, Youcef, et al.[3]The considers frequent itemsets mining in value-based databases. It shows another exact Single Scan approach for Frequent Itemset Mining (SSFIM), a heuristic as an alternative approach(EA-SSFIM), similarly as a parallel utilization on Hadoop packs (MR-SSFIM). EA-SSFIM and MR SSFIM target pitiful and huge databases, independently. The proposed philosophy (in the whole of its varieties) requires only a solitary breadth to remove the contender itemsets, and it has the favored point of view to make a fixed number of candidate itemsets openly from the estimation of the base help. This stimulates the yield system appeared differently in relation to existing approachs while overseeing insufficient and gigantic databases. Numerical results demonstrate that SSFIM outmaneuvers the stand out FIM approaches while overseeing medium and immense databases. Besides, EA-SSFIM gives near execution as SSFIM while amazingly lessening the runtime for large databases. The results in like manner reveal the pervasiveness of MR-SSFIM contemplated over the current HPC-based solutions for FIM using sparse and big databases.

According to Xun, Yaling, Jifu Zhanget.al [4] Existing parallel mining algorithms for frequent itemsets lack a mechanism that allows automatic parallelization, load effort, data distribution, and fault tolerance on large clusters. As a solution to the present disadvantage, we tend to tend to vogue a parallel frequent itemsets mining rule called FiDoop victimization the MapReduce programming model. to realize compressed storage and avoid building conditional pattern bases, FiDoop incorporates the frequent things ultrametric tree, instead of typical FP trees. In FiDoop, MapReduce jobs unit enforced to finish the mining task. at intervals the crucial third MapReduce job, the mappers severally decompose itemsets, the reducers perform combination operations by constructing very little ultrametric trees, and jointly the actual mining of those trees individually. we have a tendency to tend to implement FiDoop on our in-house Hadoop cluster. we have a tendency to tend to indicate that FiDoop on the cluster is sensitive to knowledge distribution and dimensions, as a results of itemsets with entirely completely different lengths have different decomposition and construction costs. to enhance FiDoop's performance, we tend to tend to develop a employment balance metric to measure load balance across the cluster's computing nodes. we tend to develop FiDoop-HD, Associate in Nursing extension of FiDoop, to hurry up the mining performance for high-dimensional data analysis. in depth experiments victimization real-world celestial spectral data demonstrate that our projected resolution is economical and ascendable.

Xiong, Xinyu, et al [5]Frequent itemsets mining with differential privacy refers to the substance of mining all frequent itemsets whose supports above a given threshold throughout a given transactional dataset, with the constraint that the well-mined results mustn't break the privacy of any single act. Current solutions for this disadvantage cannot well balance efficiency, privacy, information and knowledge and dat utility over large-scale information. Toward this end, we've a bent to propose a cost-effective, differential personal frequent itemsets mining rule over large-scale data. supported the concepts of sampling and dealing truncation victimization length constraints, our algorithmic program reduces the computation intensity, reduces mining sensitivity, and so improves knowledge utility given a set privacy budget. Experimental results show that our algorithmic program achieves higher performance than prior approaches on multiple dataset.

Tribhuvan et.al [6] Data mining tools estimate upcoming trends and behaviors, allowing businesses to construct practical, knowledge-driven decisions. Association Rule mining is a very important data mining practice in different fields. In most of the data mining applications, finding frequent itemsets is the crucial issue needed to be addressed. Several algorithms like Apriori, FP-Growth, and FUIT offered better solutions for mining frequent itemsets. Still, the features like automatic parallelization, fine load balancing, and distribution

of data on large clusters, are needed to be achieved. Improved Apriori Algorithm with MapReduce framework is used to solve these issues. Thus, it is possible to achieve parallelism and lessen the execution time. Here, number of mapreduce jobs are used to discover frequent itemsets of big datasets with the help of multiple computing nodes by applying parallelism among them. In this paper, the proposed system works on multiple nodes efficiently and execution time also reduced

Nikam Pallavi V., and Deepa S. Deshpande. [7] To find out the frequent itemset is very important task in data mining. These frequent itemsets are useful in applications like Association rule mining and co-relations. To extract frequent itemsets these systems are using some algorithms, but these are inefficient in distributing and balancing the load, when it comes across excessive data. Automatic parallelization is additionally unimaginable with these algorithms. to beat these problems with existing rules there's have to be compelled to develop algorithm which is able to support the missing options, like mechanically parallelization, leveling and sensible distribution of knowledge. they're employing a new approach to search out frequent itemsets by exploitation MapReduce. changed Apriori rule is employed with HDFS setting this is often known as FiDoop Technique. during this technique mapreduce method can work severally and at the same time by exploitation the decompose strategy. The result of this mapreduce technique will be given to the reducers and reducers will show the result. In the experiment they used three different algorithms like basic apriori, FP Growth and our proposed modifies apriori, the system has executed in standalone machine as well as distributed environment and shown the results how proposed algorithm is better than existing algorithms.

Zhang, and Hou Ying [8]. With the coming of big data time, the huge number of IDS log makes the traditional computing technology and systems cannot cope and deal with the needs of the analysis of security log, so large-scale computing power has become a prerequisite for the effective implementation of data mining technology. Based on the Hadoop framework, the applies the parallel frequent item sets mining algorithm to the Snort Intrusion Detection System, which solves the problem that Snort-IDS cannot judge the security event itself. At the same time, it also solves the problem of decreasing the processing speed due to the enormous increase of data. So that the system has the ability of detecting new intrusions, enrich and improve the Snort-IDS functional system and the performance.

Christian Borgelt [9] 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 oftransactions. Then select all transactions that contain the least frequent item, delete this item from them, and recourse to process the obtained reduced database, remembering that the item sets found in the recursion share the item as aprefix.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, isexploited 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 mainadvantage 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 almosttrivial 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 yieldscompetitive execution times.

Luna, José María, et al. [10] Pattern mining is one among the foremost very important tasks to extract important and useful knowledge from data. This task aims to extract item-sets that represent any form of homogeneity and regularity in information. though many economical algorithms square measure developed throughout this regard, the growing interest in knowledge has caused the performance of existing pattern mining techniques to be born. The goal of this paper is to propose new economical pattern mining algorithms to work in immense information. to the current aim, a series of algorithms supported the MapReduce framework and also the Hadoop ASCII text file implementation are planned. The planned algorithms square measure typically divided into three main groups. First, algorithms with no pruning strategy square measure planned, that extract any existing itemset in knowledge. Second, a pair of algorithms (space pruning AprioriMR and prime AprioriMR) that prune the search space by implies that of the well-known anti-monotone property ar planned. Finally, a final algorithm (maximal AprioriMR) is to boot planned for mining condensed representations of frequent patterns. to check the performance of the planned algorithms, a varied assortment of huge knowledge datasets are thought of, comprising up to transactions and over five million of distinct single-items. The experimental stage includes comparisons against very economical and well-known pattern mining algorithms. Results reveal the interest of applying MapReduce versions once advanced problems ar thought of, and put together the standard of this paradigm once managing tiny knowledge.
.

## III. Proposed System Overview

In the proposed research work to design and implement a system for FIM using privacy preservation approach. This work also carried out an efficient, differential private frequent itemsets mining algorithm over large scale data. Based on the ideas of sampling and transaction truncation using length constraints, our algorithm reduces the computation intensity, reduces mining sensitivity, and thus improves data utility given a fixed privacy budget.
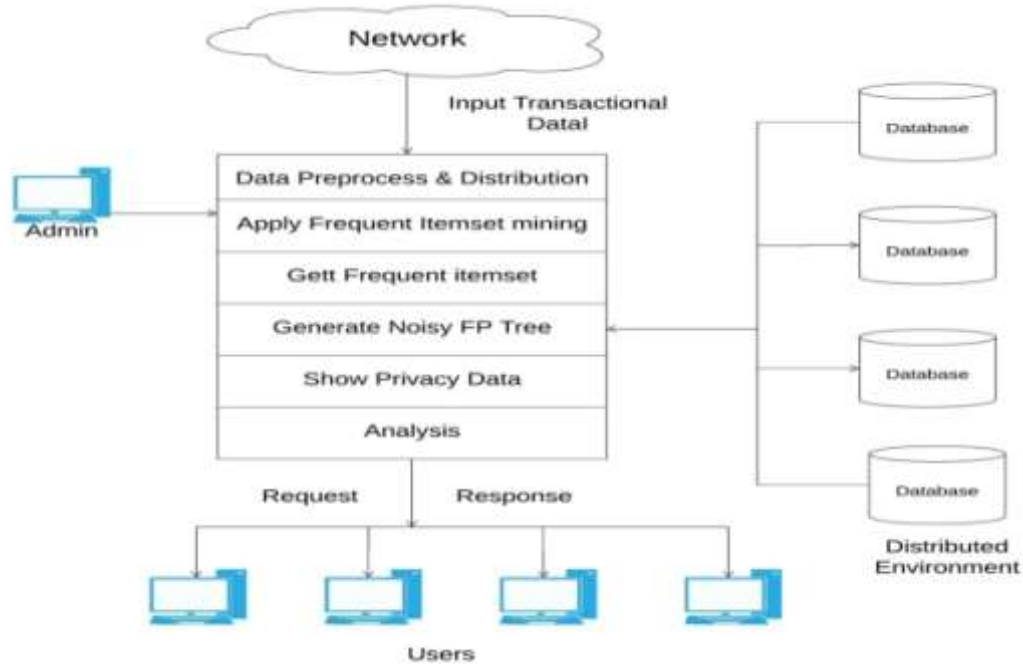


**Figure 1** : Proposed system architecture

Mining has turned into a blasting subject of research in Computer Science. This fast growing phenomenal is driven by numerous reasons. First, information mining and information warehousing are to a great degree fertile with inquire about issues but then present an outrageous helpful instrument to oversee huge measure of information. Besides, information develops at an exponential rate and Internet technology has made it simple to suspect that information from everywhere throughout the world so companies and associations these days find themselves immersed with information, and anxious to extricate useful information from them to benefit their business.

## IV. Algorithms Design

**Algorithms1 :Modified Apriori**
**Input**: Dataset D, Support generation denominator De, min_req_itemsmk;
**Output**: Generate T item set
**Step 1**: for all (T in DBi) do
**Step 2:** items [] ←split (T)
**Step 3:** for all (item in items []) do
**Step 4:** if the item is available in FLits
**Step 5:** Add a [] ← item
**Step 6:** end for
**Step 7:** add all a toArrayList<items name, count>All items.
**Step 8:** Generate the support base on support= (T.count/De)
**Step 9:** for (k in Array List)
**Step 10:** if (k.count>=support)
**Step 11:**FreqItems (k);
**Step 12**: end for
**Step 13**: apply step 9 to 12 when reach mk

# V. Results and Discussions

We used synthetic data resemble market basket data with short frequent patterns. The other two datasets are real data, which are dense in long frequent patterns. These data sets were often used in the previous study of association rules mining. The experimental results of this framework are set to the minimum support threshold (or, proportionally, bigger data sizes) than having even yet been considered. These upgrade same at no execution cost, as prove by the way that our implementation achieves the performance compare to other methods less time. Therefore we are taking proposed algorithm; it can be best algorithm to give the accurate results as compare to existing systems. Proposed system algorithm shows the faster execution even for large database. We could create our own vast dataset against which to likewise run tests; however the cost for doing as such is negligible. The information in the web documents set comes from a real domain and so is meaningful. We have used five sets of data in our experiments. Three of these sets are synthetic data T10I4D100K, T25I10D10K, and T40I10D100k. The other two datasets are real data (Groceries) which are dense in long frequent patterns. These data sets were live data sets for the study of association rules mining and were downloaded from http://www.jbtraders.in/ .

For the proposed execution we have used three different dataset, the grocery dataset has taken from www.jbtraders.in and some electronic item base synthetic dataset has given from internet. The third dataset has taken from sport www.sports365.in. Below we have done multiple experiments which are represented in graphs.
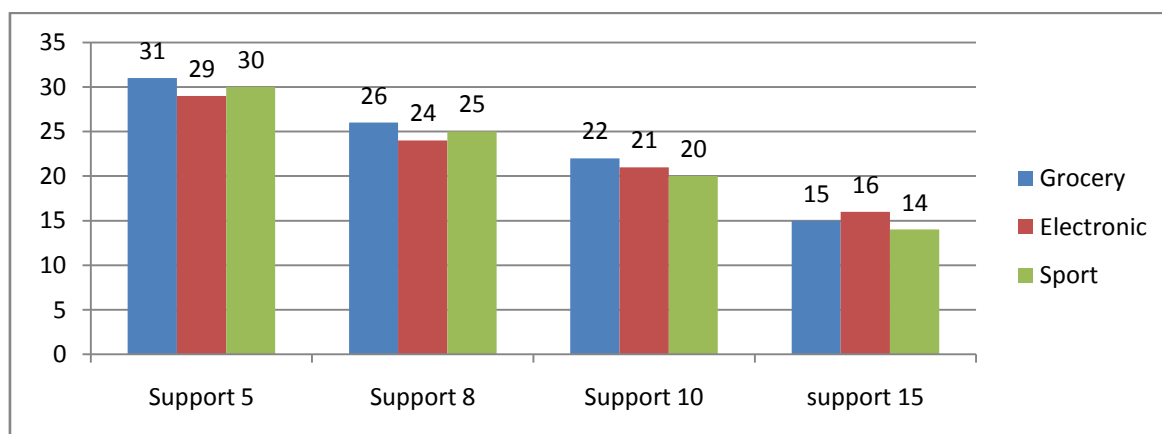
**Performance of Fidoop in terms of Time required in seconds with different support denominator with different dataset**

Following table indicates the time taken in seconds by Proposed Algorithm proposed for the three different datasets of Grocery, Electronic & Sports for different support counts.

**Table 1: Performance of Time required in seconds with different support denominator with different dataset**

| Support Count | Time in Seconds | | Time in Seconds | | Time in Seconds | |
|---|---|---|---|---|---|---|
| Support 5 | 3 | 1 | 2 | 9 | 3 | 0 |
| Support 8 | 2 | 6 | 2 | 4 | 2 | 5 |
| Support 10 | 2 | 2 | 2 | 1 | 2 | 0 |
| support 15 | 1 | 5 | 1 | 6 | 1 | 4 |

Figure2 shows the result presented in Table 1 for three different item set of Grocery, Electronic & sports for various support values of 5,8,10 and 15% respectively. It shows the time required to extract the frequent item set in seconds with different dataset. Utilization is very good in applying the algorithm Fidoop for frequent item set mining with retail item set for various support values.



**Figure 2: Time required in seconds with different support denominator with different dataset.**

The experiment results showed using graphs exhibits that the time used for extraction of frequent item set using Fidoop decreases as the support increases.

## VI. Conclusion

In this work, system investigates the problem of designing a differentially private FIM algorithm. We use differential privacy to stop the potential information exposure about individual record set during the data mining process. Here we studied system model of Frequent Item set Mining using distributed environment. We put forward algorithm and mining long patterns.. It minimizes time required for high dimensional dataset. As we are using map reduce here, can also handle huge size dataset without any problem. We represented comparative table between different algorithms used in FIM. As our future work we plan to design more effective differentially private FIM on big data.

## References

[1]. Pan, Zhaopeng, Peiyu Liu, and Jing Yi. "An Improved FP-Tree Algorithm for Mining Maximal Frequent Patterns."2018 10th International Conference on Measuring Technology and Mechatronics Automation (ICMTMA).IEEE, 2018.

[2]. Fahrudin, TresnaMaulana, IwanSyarif, and Ali RidhoBarakbah. "Discovering patterns of NED-breast cancer based on association rules using apriori and FP-growth." Knowledge Creation and Intelligent Computing (IES-KCIC), 2017 International Electronics Symposium on.IEEE, 2017.

[3]. Djenouri, Youcef, et al. "Frequent Itemset Mining in Big Data With Effective Single Scan Algorithms." IEEE Access 6 (2018): 68013-68026.

[4]. Xun, Yaling, Jifu Zhang, and Xiao Qin. "Fidoop: Parallel mining of frequent itemsets using mapreduce." IEEE transactions on Systems, Man, and Cybernetics: systems 46.3 (2016): 313-325.

[5]. Xiong, Xinyu, et al."FrequentItemsets Mining with Differential Privacy over Large-scale Data." IEEE Access (2018).

[6]. Nikam, Pallavi V., and Deepa S. Deshpande. "New Approach in Big Data Mining for Frequent Itemset Using Mapreduce in HDFS."2018 3rd International Conference for Convergence in Technology (I2CT).IEEE, 2018.

[7]. Tribhuvan, Seema A., Nitin R. Gavai, and Bharti P. Vasgi. "Frequent Itemset Mining Using Improved Apriori Algorithm with MapReduce." 2017 International Conference on Computing, Communication, Control and Automation (ICCUBEA).IEEE, 2017.

[8]. Tong, Zhang, and Hou Ying. "Application of frequent item set mining algorithm in IDS based on Hadoop framework." 2018 Chinese Control And Decision Conference (CCDC). IEEE, 2018.

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

[10]. Luna, José María, et al. "Apriori versions based on mapreduce for mining frequent patterns on big data." IEEE transactions on cybernetics (2017).