

Frequent Item sets Mining with Differential Privacy over Large Scale Data

¹PRAVALLIKA BUDAVARAPU, ²A.BINDUKALA, ³A.D.SIVARAMA KUMAR

¹M.Tech Student, ^{2,3}Assistant Professor
DEPT OF CSE
SVR Engineering College, Nandyal

Abstract—Frequent itemsets mining with differential privacy refers to the problem of mining all frequent itemsets whose supports are above a given threshold in a given transactional dataset, with the constraint that the mined results should not break the privacy of any single transaction. Current solutions for this problem cannot well balance efficiency, privacy and data utility over large scaled data. Toward this end, we propose 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. Experimental results show that our algorithm achieves better performance than prior approaches on multiple datasets.

Index Terms—Frequent Itemsets Mining; Differential Privacy; Sampling; Transaction Truncation; String Matching

I.INTRODUCTION

In recent years, with the explosive growth of data and the rapid development of information technology, various industries have accumulated large amounts of data through various channels. To discover useful knowledge from large amounts of data for upper-layer applications (e.g .business decisions, potential

customer analysis, etc.), data mining has been developed rapidly. It has produced a positive impact in many areas such as business and medical care. Along with the great benefits of these advances, the large amount of data also contains privacy sensitive information, which may be leaked if not well managed. For instance, smart phone applications are recording the whereabouts of users through GPS sensors and are transferring the data to their servers. Medical records are also storing potential relationships between diseases and a variety of data. Mining on user location data or medical record data both provide invaluable information; however, they may also leak user privacy. Thus mining knowledge under confident privacy guarantees is highly expected. This s investigates how to mine frequent item sets with privacy guarantee for big data. I consider the following application scenario. A company (such as information consulting firm) has a large-scale dataset. The company would like to make the dataset public and therefore allow the public to execute frequent item sets mining for getting cooperation or profits. But due to privacy considerations, the company cannot provide the original dataset directly. Therefore, privacy mechanisms are needed to process the data, which is the focus of this project. To ensure privacy of data mining, traditional methods are based on k-anonymity and its extended models. These methods require certain assumptions; it is difficult to protect

privacy when the assumptions are violated. The insufficiency of k-anonymity and its extended models is that there is no strict definition of the attack model, and that the knowledge of the attacker cannot be quantitatively defined. To pursue strict privacy analysis, work proposed a strong privacy protection model called differential privacy. This privacy definition features independence of background knowledge of the attacker and proves very useful. Frequent pattern mining with privacy protection has also received extensive attention. As preliminary methods, these works have provided a lot of contributions in this area. But with the advance of research, these privacy method shave not been able to provide effective privacy. In order to overcome these difficulties, researches began to focus on the differential privacy protection framework. Although guaranteeing privacy temporary, however, the balance between privacy and utility of frequent item sets mining results needs to be further pursued.

In this project, I propose a novel differential private frequent item sets mining algorithm for big data by merging the ideas of, which has better performance due to the new sampling and better truncation techniques. Build algorithm on FP-Tree for frequent item sets mining. In order to solve the problem of building FP-Tree with large-scale data, I first use the sampling idea to obtain representative data to mine potential closed frequent item sets, which are later used to find the final frequent items in the large-scale data. In addition, I employ the length constraint strategy to solve the problem of high global sensitivity. Specifically, I use string matching ideas to discover the most similar string in the source dataset, and implement transaction truncation for achieving the lowest information loss. I finally add the Laplace

noise for frequent item sets to ensure privacy guarantees.

A few challenges exist: First, how to design a sampling method to control the sampling error? I use the central limit theorem to calculate a reasonable sample size to control the error range. After obtaining the sample size, the dataset is randomly sampled using a data analysis toolkit. The second challenge is how to design a good string matching method to truncate the transaction without losing information as far as possible? I match the potential item sets in the sample data to find the most similar items and then merge them with the most frequent items until the maximum length constraint is reached.

EXISTING SYSTEM:

Explosive growth of data and the rapid development of information technology, various industries have accumulated large amounts of data through various channels. To discover useful knowledge from large amounts of data for upper-layer applications (e.g. business decisions, potential customer analysis, etc.), data mining has been developed rapidly.

It has produced a positive impact in many areas such as business and medical care. Along with the great benefits of these advances, the large amount of data also contains privacy sensitive information, which may be leaked if not well managed.

The company would like to make the dataset public and therefore allow the public to execute frequent item sets mining for getting cooperation or profits.

DISADVANTAGES:

The system doesn't work efficiently.

There is no an affective privacy preserving encryption techniques in this system.

III. PROPOSED SYSTEM:

I propose a novel differential private frequent item sets mining algorithm for big data by merging the ideas, which has better performance due to the new sampling and better truncation techniques.

Build algorithm on FP-Tree for frequent item sets mining. In order to solve the problem of building FP-Tree with large-scale data, I first use the sampling idea to obtain representative data to mine potential closed frequent item sets, which are later used to find the final frequent items in the large-scale data

ADVANTAGES OF PROPOSED SYSTEM:

Frequent Item Sets Mining With Differential Privacy Over large Scale Data which can the collusion attack launched by the users.

Attackers can only observe the encrypted data stored in the cloud. In order to avoid, the well-known cipher text-only attack model has been implemented.

IV.SYSTEM DESIGN

4.1 System Architecture:

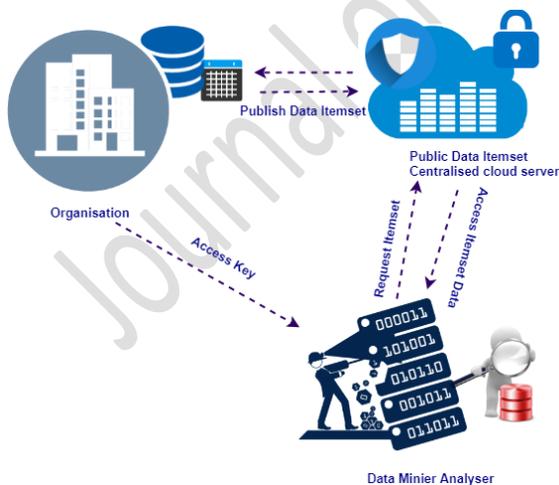


Fig : 4.1 System Architecture

V. MODULES

Admin

In this module, the Cloud has to login by using valid user name and password. After login successful he can do some operations such as List all users and authorize, View all company users and authorize Add all company name and view, View all company details with rank and reviews, View all companies by Frequent Item sets Mining using FP-Tree format and give link on company name view its details, View all user search transaction by keyword, Show search ratio by keyword, Find top k Frequent item sets by ranks View all companies rank by chart, View all search ratio by keyword in chart

Production Company

In this module, there are n numbers of Owners are present. Owner should register before doing any operations. Once registers, their details will be stored to the database. After registration successful, he has to login by using authorized user name and password. Once Login is successful Owner will do some operations like View your profile, Add company data set, View your company details with reviews and rank, View user search transactions on your company, View other related companies by Frequent Item sets Mining using FP-Tree format and give link on company name view its details

Users

In this module, there are n numbers of users are present. User should register before doing any operations. Once user registers, their details will be stored to the database. After registration successful, he has to login by using authorized user name and password. Once Login is successful user will do some

operations like View your profile, Search companies by keyword and show all related companies by FP-Tree format and give link on company name view its details, view its details with image(increment rank),review , show other review also, find search ratio, View your search transactions by keyword

VI.CONCLUSION

In this project, I propose a novel differentially private algorithm for frequent item sets mining. The algorithm features better data utility and better computation efficiency. Various experimental evaluations validate that the proposed algorithm has high F-Score and low relative error. A lesson learned is that fine tuned parameters lead to better differentially private frequent item sets mining algorithms with regard to data utility.

REFERENCES

- [1] Z. John Lu, "The elements of statistical learning: data mining, inference, and prediction," *Journal of the Royal Statistical Society: Series A(Statistics in Society)*, vol. 173, no. 3, pp. 693–694, 2010.
- [2] U. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, "From data mining to knowledge discovery in databases," *AI magazine*, vol. 17, no. 3, p. 37, 1996.
- [3] H. Yang, K. Huang, I. King, and M. R. Lyu, "Localized support vector regression for time series prediction," *Neurocomputing*, vol. 72, no. 10-12, pp. 2659–2669, 2009.
- [4] C. Romero and S. Ventura, "Educational data mining: A review of the state of the art," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 40, pp. 601–618, Nov 2010.
- [5] J. Han, J. Pei, and M. Kamber, *Data mining: concepts and techniques*. Elsevier, 2011.

- [6] X. Fang, Y. Xu, X. Li, Z. Lai, and W. K. Wong, "Robust semi-supervised subspace clustering via non-negative low-rank representation," *IEEE Transactions on Cybernetics*, vol. 46, pp. 1828–1838, Aug 2016.
- [7] M. Peñna, F. Biscarri, J. I. Guerrero, I. Monedero, and C. Le'on, "Rule based system to detect energy efficiency anomalies in smart buildings, a data mining approach," *Expert Systems with Applications*, vol. 56, pp. 242–255, 2016.
- [8] Y. Guo, F. Wang, B. Chen, and J. Xin, "Robust echo state networks based on correlation entropy induced loss function," *Neurocomputing*, vol. 267, pp. 295–303, 2017.
- [9] H. Lim and H.-J. Kim, "Item recommendation using tag emotion in social cataloging services," *Expert Systems with Applications*, vol. 89, pp. 179–187, 2017.
- [10] R. C.-W. Wong, J. Li, A. W.-C. Fu, and K. Wang, "(ϵ , k)-anonymity: an enhanced k-anonymity model for privacy preserving data publishing," in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 754–759, ACM, 2006.