

# Privacy Preserving in Data Sets Through Multiple Shuffle

S. Gayathri<sup>1</sup>, P. Rakini Prasadha<sup>2</sup>, G. Vinitha<sup>3</sup>

<sup>1</sup>Assistant Professor, Dept. of Computer Science and Engg., Sri Eshwar College of Engg., Coimbatore, India

<sup>2,3</sup>U.G. Student, Dept. of Computer Science and Engg., Sri Eshwar College of Engg., Coimbatore, India

**Abstract:** Cloud computing and its model for IT services supported the net and large knowledge centers, the outsourcing of information and computing services. A company (data owner) lacking in expertise or machine resources will source its mining has to a 3rd party service supplier (server). However, each the things and also the association rules of the outsourced information area unit thought of holding of the corporation (data owner). Big knowledge privacy-preserving has attracted increasing attention of researchers in recent years. But existing models area unit thus difficult and long that they're challenging to implement. In this paper, we propose a more feasible and efficient model for big data sets privacy-preserving using shuffling multiple attributes (M-Shuffle) to achieve a tradeoff between data utility and privacy. Our strategy is, first of all, categorize all the records into some teams victimization K-means formula per the sensitive attributes. Then we choose the columns to be shuffled using entropy. At last, we introduce the random shuffle algorithm to our model to break the correlation among the columns of big data sets. To protect company privacy, information owner transforms its data and ships it to the server, sends mining queries to the server, and recovers verity patterns from the extracted patterns received from the server. In this paper, we have a tendency to study the matter of outsourcing the association rule mining task among a company privacy-preserving framework. We propose AN attack model supported information and devise a theme for privacy conserving outsourced mining. Our theme ensures that every reworked item is indistinguishable with regard to the attacker's information, from a minimum of  $k-1$  alternative reworked things.

**Keywords:** data sets

## 1. Introduction

In the information age, big data is a milestone and leads to sharp changes in modern society. Government agencies, big IT companies, and other organizations always publish big data sets for research purpose, for example, the census or medical datasets. But releasing the datasets to the public may cause privacy leakage because every record stored in such kind of data sets corresponds to one specific individual. Privacy-preserving in big data sets is hence to become a big challenge worldwide.

There are two main branches in privacy research: privacy-preserving data publishing (PPDP) and privacy-preserving data mining (PPDM). In PPDP, there are several milestones like  $k$ -anonymity,  $l$ -diversity,  $t$ -closeness and differential privacy. The extensions of this models are also widely studied in recent years. The  $k$ -anonymity and its extensions are likely to suffer

homogeneity attack and background knowledge attack, which makes it very vulnerable. The  $l$ -diversity and its extensions sometimes lose more information and lead to a larger utility loss. Although  $t$ -closeness offers better privacy protection, it suffers the same challenge of  $l$ -diversity on some occasions. Differential privacy offers a theoretical foundation but sometimes it is too strict to implement on real-world systems and the efficiency is another big challenge.

However, there are three urgent challenges needs to be solved. The first challenge is that privacy and data utility seems to be a natural antithesis. A tradeoff must be found to preserve the privacy of big data sets while the utility of the big data sets should be maintained at a proper level. It's hard to satisfy all the requirements.

The second challenge is that the correlation between the values of the same record, anonymization, and generalization could partly decouple the correlation between the values, but the utility will suffer a great loss in this. The time complexities of this kind of methods are very large, which is a big problem for practical usage.

The third challenge occurs when adversaries have too much background knowledge. Background knowledge is a great threat to all kinds of privacy models because an adversary can re-identify a specific person if he gains some information from the released datasets and combines it with the background knowledge of himself.

In order to address these challenges, we propose a multiple shuffle model. In the M-Shuffle, we first use  $K$ -means to group all the records to  $K$  clusters. Then, we group all the attributes and choose some of the attributes to be shuffled using entropy. At last, we shuffle the chosen columns using random shuffle algorithm. By applying the mechanism, we break the correlation among the values of one single record. In this way, background knowledge becomes less useful and the statistics of a single column will not change at all. Using  $K$ -means will improve the utility without sacrificing the privacy level. Our model will maintain the statistics of the shuffled columns and provide the desired privacy level with high efficiency.

## 2. Related work

With the development of information technology, the big data age has arrived. Its impacts are so pervasive that we can

see its implementations on every aspect of daily life, research, or even government functioning. For example, existing big data sets benefit us a lot in biology, social science, e-commerce, disease-control, and so on. We can easily predict the outburst of an infectious disease through social networks which are big data sets in essence. In 2011, Doug Laney proposed an early concept of big data in the Gartner report, where big data was defined as large and complex data sets that current computing facilities were not able to handle. This started the research enthusiasm for big data. Big data never asks why while gives the predictions simply, which makes research on it more valuable.

But releasing the big data sets to the public may lead to privacy disclosure, even for research purpose. There are two main kinds of information disclosure: identity disclosure and attributes disclosure. No matter which of these two situations occurs, it will be harmful to an individual's privacy and may even cause financial loss. Privacy study has sprung up since two decades ago. First, data clustering methods were carried out to privacy-preserving. The first milestone is the  $k$ -anonymity model. It was the first model to introduce the data clustering method to privacy protection, which was proposed in 1998. The following one is  $L$ -diversity, which showed up in 2007. It is an extension of  $k$ -anonymity and introduced diversity into data clustering. Then  $t$ -closeness was employed in 2010 which also took distribution into consideration. Models based on data clustering advanced the privacy-level of big datasets. Despite their feasibility, lack of firm theoretical analysis is always a flaw.

In 2006, another milestone with the firm foundation was developed. Dwork proposed differential privacy.

Differential privacy may be a framework for formalizing privacy in applied math databases introduced so as to safeguard against these types of deanonymization techniques.

After this, extensions like personalized differential privacy frameworks appeared.

### 3. System modeling and analysis

#### A. Dataset Collection and Encryption

A knowledge set (or data set) could be an assortment of knowledge. Dataset is collected from Belgium retail market dataset. It contains the (anonymized) retail market basket knowledge from associate anonymous Belgian business establishment.

The data are provided 'as is'. Basically, any use of the information is allowed as long because the correct acknowledgment is provided and a replica of the work is provided to Tom Brijs.

The grocery store carries sixteen,470 distinctive SKU's, however a number of them solely on a seasonal basis. In total, 5,133 customers have purchased a minimum of one product within the grocery store throughout the information assortment amount.

A dataset is encrypted by using Homomorphic encryption.

Homomorphic encryption is an encryption scheme which transforms a TDB  $D$  into its encrypted version  $D^*$ .

#### B. Grouping Items for $k$ -Privacy

Given the things support table, many ways may be adopted to cluster the things into teams of size  $k$ .

We start from a simple grouping method. We assume the item support table is sorted in descending order of support and refer to cipher items in this order as  $e_1, e_2$ , etc.

Assume  $e_1, e_2 \dots e_n$  is the list of cipher items in descending order of support (with respect to  $D$ ), the groups created are, and so on. The last group, if less than  $k$  size is merged with its previous group.

Given the fact that the support of the items strictly decreases monotonically, the grouping is optimal among all the groupings with the item support table sorted in descending order of support. This means, it minimizes  $\|G\|$ , the size of the fake transactions added, and hence the size  $\|D^*\|$ .

#### C. Constructing Fake Transactions

Given a noise table specifying the noise  $N(e)$  needed for each cipher item  $e$ , we generate the fake transactions as follows.

First, we have a tendency to drop the rows with zero noise, corresponding to the most frequent items of each group or to other items with support equal to the maximum support of a group.

Second, we have a tendency to type the remaining rows in descending order of noise.

This technique yields a minimum variety of various varieties of pretending transactions that equal the number of cipher things with distinct noise. This observation yields a compact abstract for the shopper of the introduced pretend transactions.

The purpose of employing a compact abstract is to cut back the storage overhead at the facet of the information owner UN agency might not be equipped with ample machine resources and storage, that is common within the outsourcing information model.

In order to implement the synopsis efficiently, we use a hash table generated with a minimal perfect hash function.

Minimal good hash functions are widely used for memory economical storage and quick retrieval of things from static sets.

A marginal good hash perform could be a good hash perform that maps  $n$  keys to  $n$  consecutive integers, usually  $[0 \dots n - 1]$ .

#### D. Decryption

When the shopper requests the execution of a pattern mining question to the server, specifying a minimum support threshold  $\sigma$ , the server returns the computed frequent patterns from  $D^*$ .

Clearly, for every item set  $S$  and its corresponding cipher item- set  $E$ , we have that  $\text{supply}(S) \leq \text{supply}^*(E)$ .

For each cipher pattern  $E$  returned by the server together with  $\text{supply}^*(E)$ , the  $E/D$  module recovers the corresponding plain pattern  $S$ .

It has to reconstruct the precise support of  $S$  in  $D$  and choose

on this basis if  $S$  could be a frequent pattern.

To achieve this goal, the E/D module adjusts the support of  $E$  by removing the effect of the fake transactions.

#### 4. System analysis

##### A. Data Clustering Algorithm

K-means is a very useful tool of data clustering; it is first used in signal processing. K-means aims to partition  $n$  records into  $k$  clusters in which each record belongs to the cluster with the nearest mean. The iterative update is used in the most common algorithm. The algorithm, k-means, is named because of its extensive existence. It is also referred to Lloyd's algorithm, particularly in the computer science community.

##### B. Shuffle Algorithm

The Fisher-Yates shuffle is a very popular algorithm which is an in-place shuffle. That means instead of creating a new shuffled copy of the records, it shuffles the records of a table in place. If the table to be shuffled is large enough, this mechanism can fit well.

In order to initialize and shuffle a table synchronously, and the advanced version is introduced to our mechanism to make it more efficient. The random algorithm can perfectly put a certain record  $i$  into a random location among the first  $i$  locations in the table, after moving the record previously taking up that location to location  $i$ . In normal conditions, which records to be shuffled by a column of number, especially the integers, this could be easy to be represented by a function because the implementation will not change it.

#### 5. Conclusions and future work

We present that existing models are complex and hard to implement. Therefore, we propose a more feasible and practical model using multiple shuffling. In this paper, we introduce M-Shuffle. This model decouples the correlation among values of a column so that we can protect privacy without damaging the statistics. Experiments on real-world datasets show the effectiveness and efficiency of the proposed model.

The project involved the problem of (corporate) privacy-preserving mining of frequent patterns (from which association rules can easily be computed) on an encrypted outsourced TDB. We assumed that a conservative model wherever someone is aware of the domain things of things and their actual frequency and may use this data to spot cipher items and cipher item sets.

We proposed an encryption scheme, called Homomorphic, which is based on 1-1 substitution ciphers for items and adding fake transactions to make each cipher item share the same frequency.

It makes use of a compact precis of the pretend transactions from that truth support of strip-mined patterns from the server will be expeditiously recovered.

We also proposed a strategy for incremental maintenance of the synopsis against updates consisting of appends and dropping off old transaction batches.

Currently, our privacy analysis relies on the idea of the equal probability of candidates.

It would be fascinating to boost the framework and also the analysis by appealing to scientific discipline notions like good secrecy.

Moreover, our work considers the ciphertext-only attack model, within which the offender has access solely to the encrypted things. It may well be fascinating to contemplate alternative attack models wherever the wrongdoer is aware of some pairs of things and their cipher values.

There are several directions for future work. The first one is that we want to carry out a better shuffle algorithm so that the model can reach a better privacy level without sacrificing utility. Then the measurement of utility is a big challenge to modern privacy study. That's why we want to propose a universe measurement by combining some modern theories, for example, the game theory, the information theory and so on.

#### References

- [1] L. Xu, C. Jiang, J. Wang, J. Yuan and Y. Ren, "Information Security in Big Data: Privacy and Data Mining," in *IEEE Access*, vol. 2, pp. 1149-1176, 2014.
- [2] P. Samarati and L. Sweeney. Protecting privacy once revealing information: k-anonymity and its social control through generalization and suppression, In *Proceedings of the IEEE Symposium on Research in Security and Privacy*, pages 1-19, 1998.
- [3] A. Machanavajhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian. L-diversity: Privacy beyond k-anonymity. *ACM Trans. Knowl. Discov. Data.*, 1(1), 2007.
- [4] N. Li, T. Li and S. Venkatasubramanian. "Closeness: A New Privacy Measure for Data Publishing," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 7, pp. 943-956, July 2010.
- [5] Shui Yu. Big privacy: Challenges and opportunities of private study within the age of massive knowledge. *IEEE Access*, 4:2751-2763, June 2016.
- [6] V. Marx. Biology: The big challenges of big data. *Nature*, 498:255-260, 2013.
- [7] G. King. Ensuring the data-rich future of the social sciences. *Science*, 331(6018):719-721, 2011.