

Hiding Sensitive Patterns in Association Rules Mining

Guanling Lee*

Chien-Yu Chang*

Arbee L.P Chen⁺

*Department of Computer Science
and Information Engineering
National Dong Hwa University

Email: guanling@mail.ndhu.edu.tw

⁺Department of computer
Science

National Cheng-chi University

Email: alpchen@cs.nthu.edu.tw

Abstract

Data mining techniques have been developed in many applications. However, it also causes a threat to privacy. We investigate to find an appropriate balance between a need for privacy and information discovery on association patterns. In this paper, we propose an innovative technique for hiding sensitive patterns. In our approach, a sanitization matrix is defined. By multiplying the original transaction database and the sanitization matrix, a new database, which is sanitized for privacy concern, is gotten. Moreover, a set of experiments is performed to show the effectiveness of our approach.

Keywords: association patterns, privacy preservation, sanitized database, data mining.

1. Introduction

Data mining techniques have been developed in many applications and researches. However, it also brings the problem of privacy. A motivating example is discussed in [4]. Suppose we have a server and many clients in which each client has a set of data. The clients want the server to gather statistical information about association among items in order to provide recommendations to the customers. However, the clients do not want the server to know some *sensitive patterns*. Sensitive pattern is the

frequent itemset that contain highly sensitive knowledge. Thus, when a client sends its database to the server, some sensitive patterns are hidden from its database according to some specific privacy policies. Therefore, the server only can gather statistical information from the modified database.

In recent years, more and more researchers emphasize the seriousness of the problem about privacy. The privacy problem can be classified into two classes: *data privacy* problem and *information privacy* problem. Data privacy is to protect the privacy of sensitive data, while information privacy is investigated privacy of patterns that contain highly sensitive knowledge.

Privacy-preserving mining in the context of data privacy for classification rules has been investigated in [3]. By using a randomizing function with Gaussian or Uniform perturbations, the sensitive values in user's record will be perturbed. Based on probabilistic distortion of user data, [6] demonstrates a scheme. In [5], the problem of how to avoid privacy breaches in privacy preserving data mining is introduced.

Information privacy preserving problem is to hide the sensitive patterns or rules by updating the original database and with as little effect on non-sensitive patterns as possible. This problem is

proved to be NP-Hard [ABE99]. Similar to [ABE99], the other heuristic method is proposed in [8]. They falsify some value or replace known values with unknown values such as question marks.

In [7], a framework is proposed to enforce privacy in mining frequent itemsets. In the approach, the victim items that should be eliminated for each restrictive pattern are selected. And transaction retrieval engine is used to identify sensitive transactions for each restrictive pattern. Based on the disclosure threshold, the number of sensitive transactions is computed and the victim items are removed from the select transactions. In this paper, we propose an innovative technique for hiding sensitive patterns. By observing the relationship between sensitive patterns and non-sensitive patterns, a *sanitization matrix* is defined. By setting the entries in sanitization matrix to appropriate values and multiplying the original transaction database to the sanitization matrix, a *sanitized database* is gotten. The sanitized database is the database that has been modified for hiding sensitive patterns with privacy concern.

The reminder of this paper is organized as follows. The problem and the framework of our

approach is presented in section 2. In section 3, the sanitizing algorithms are discussed. The metrics to estimate the performance of our approach is introduced in section 4. The experimental results are also reported in section 4. We conclude with a summary and directions for future work in section 5.

2. Basic Concept

2.1 Problem Formulation

In our approach, a transaction database \mathbf{D} is represented as a matrix in which the rows represent transactions and the columns represent the items. If \mathbf{D} contains m transactions and n kinds of items, \mathbf{D} is represented as an $m \times n$ matrix. The entry $D_{t,i}$ is set to 1 if item i is purchased in transaction t . Otherwise, it is set to 0.

Our problem can be formulated as follows. Let \mathbf{D} be a transaction database, \mathbf{P} be the set of frequent patterns that can be mined from \mathbf{D} . Let \mathbf{P}_h denote a set of sensitive patterns that need to be hidden according to some security policies, and $\mathbf{P}_h \subset \mathbf{P}$. $\sim\mathbf{P}_h$ is the set of non-sensitive patterns. $\sim\mathbf{P}_h \cup \mathbf{P}_h = \mathbf{P}$. Our problem is to transform \mathbf{D} into a sanitized database \mathbf{D}' such that only the patterns belong to $\sim\mathbf{P}_h$ can be mined from \mathbf{D}' .

$$\begin{array}{c}
 \begin{array}{ccc} & 1 & 2 & 3 \\ \begin{array}{l} t1 \\ t2 \\ t3 \\ t4 \end{array} & \left(\begin{array}{ccc} 1 & 0 & 1 \\ \mathbf{1} & \mathbf{1} & 0 \\ 0 & 0 & 1 \\ \mathbf{1} & \mathbf{1} & 1 \end{array} \right) \\ \mathbf{D} & & &
 \end{array}
 \end{array}
 \times
 \begin{array}{c}
 \begin{array}{ccc} & 1 & 2 & 3 \\ \begin{array}{l} 1 \\ 2 \\ 3 \end{array} & \left(\begin{array}{ccc} 1 & 0 & 0 \\ \mathbf{-1} & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) \\ \mathbf{CM} & & &
 \end{array}
 \end{array}
 =
 \begin{array}{c}
 \begin{array}{ccc} & 1 & 2 & 3 \\ \begin{array}{l} t1 \\ t2 \\ t3 \\ t4 \end{array} & \left(\begin{array}{ccc} 1 & 0 & 1 \\ \mathbf{0} & \mathbf{1} & 0 \\ 0 & 0 & 1 \\ \mathbf{0} & \mathbf{1} & 1 \end{array} \right) \\ \mathbf{D}' & & &
 \end{array}
 \end{array}$$

Figure 1. Setting $r_{21}=-1$ in S

2.2 Sanitization matrix

In our approach, original database \mathbf{D} is multiplied by a sanitization matrix (\mathbf{S}) to get a sanitized database \mathbf{D}' . By setting S_{ij} where $i \neq j$ to appropriate value, a sanitized database \mathbf{D}' will be gotten. In the following, the basic concept of our approach is discussed.

2.2.1 New definition for the matrix multiplication

In our approach, the matrix multiplication method is defined as follows:

1. If D_{ii} equals zero, no multiplication proceeds on it. That is, D'_{ii} is set to 0 directly. This is because our goal is to hide the sensitive pattern by decreasing its support. Moreover, if an entry with value zero can be converted to 1, new patterns may be produced.
2. If the resulting value larger than 1, set it to 1.
3. If the resulting value smaller than 0, set it to 0.

2.2.2 The Setting of “-1”

A sensitive pattern can be hidden by decreasing its support. If D_{ii} and D_{ij} are both equal to 1, set D_{ii} or D_{ij} be 0 can reduce the support of $\{i, j\}$. Refer to Figure 1. Let minimum support be 50% and $\{1, 2\}$ be a sensitive pattern. If S_{21} is set to -1, D'_{21} , D'_{41} will become 0. Oppositely, if S_{12} is set to -1, D'_{22} , D'_{42} will become 0. Therefore, the support of $\{1, 2\}$ can be decreased by setting S_{21} or S_{12} to -1. Moreover, if S_{ij} is set to -1, then for a transaction t , where D_{ii} and D_{ij} are both equal to 1, D'_{ij} will be 0.

2.2.3 The Setting of “1”

Setting appropriate entries in \mathbf{S} to -1 can reduce the support of the sensitive patterns. However, it also leads to accidentally conceal the non-sensitive patterns. We remedy this defect by setting some entries in \mathbf{S} to 1 to minimize the effect on losing non-sensitive patterns. Continue above example, the frequent patterns in \mathbf{D} are $\{1, 2\}$ and $\{1, 3\}$. Let $\{1, 2\}$ and $\{1, 3\}$ be the sensitive and non-sensitive pattern, respectively. If S_{21} is set to -1, D'_{21} and D'_{41} will be 0. As a result, the support of $\{1, 3\}$ will be decreased to 25% in \mathbf{D}' and no longer be a frequent pattern. To reserve the non-sensitive pattern $\{1, 3\}$ in \mathbf{D}' , S_{31} is set to 1 to make D'_{t1} keep the same value as D_{t1} for those transaction t where $D_{t1}=1$ and $D_{t3}=1$. The purpose of setting the specific entry to 1 is to reinforce the relation of $\{1, 3\}$ and avoid eliminating $\{1, 3\}$ accidentally. Setting corresponding entries between any two items contained in non-sensitive patterns in \mathbf{S} to 1 can preserve non-sensitive patterns after the sanitization process.

3. The Sanitization Algorithms

3.1 Hidden-First Algorithm

The main idea of Hidden-First algorithm, denoted by HF, is to eliminate all patterns in \mathbf{P}_h from \mathbf{D} by setting proper entries in \mathbf{S} to -1. The entries of \mathbf{S} are set to the proper values according to the following rules.

1. $S_{ii}=1$, diagonal entry.
2. $S_{ij} = -1$, If $\exists \rho \in \mathbf{P}_h$, such that $\{i, j\} \subseteq \rho$ and $\forall \rho' \in \sim \mathbf{P}_h$, $\{i, j\} \not\subseteq \rho'$. Moreover, the number of patterns containing $\{j\}$ in $\sim \mathbf{P}_h$ is smaller than

that of the patterns containing $\{i\}$ in $\sim\mathbf{P}_h$. The reason is that by setting S_{ij} to -1 , the support of item j will be reduced. Moreover, item j has smaller effect on $\sim\mathbf{P}_h$ than item i .

3. $S_{ij} = 0$, otherwise.

Hidden-First Algorithm

Input: $\mathbf{P}_h, \sim\mathbf{P}_h, \mathbf{D}, \mathbf{S}$

Output: \mathbf{D}'

Step 1: Set the values of the entries in \mathbf{S} according to the rules.

Step 2: (matrices multiplication)

For every transaction i in \mathbf{D} do

For $j=1$ to number of items do

If $(D_{ij}=0)$ $D'_{ij}=0$;

Else

$$D'_{ij} = \max\left(\sum_{k=1}^{\text{number of items}} D_{ik} \times S_{kj}, 0\right)$$

However, in this approach, some non-sensitive patterns may be accidentally hidden due to setting the value of some entries in \mathbf{S} to -1 .

3.2 The Non-Hidden-First Algorithm (NHF)

The main idea behind the Non-Hidden-First algorithm, denoted by NHF, is to reserve all non-sensitive patterns and endeavor to hide sensitive patterns from \mathbf{D} at the same time.

The entries in \mathbf{S} are set according to the following rules.

1. $S_{ii} = 1$, diagonal entry.
2. $S_{ij} = -1$, If $\exists \rho \in \mathbf{P}_h$, such that $\{i, j\} \subseteq \rho$ and $\forall \rho' \in \sim\mathbf{P}_h$, $\{i, j\} \not\subseteq \rho'$. Moreover, the number of patterns containing $\{j\}$ in $\sim\mathbf{P}_h$ is smaller than that of the

patterns containing $\{i\}$ in $\sim\mathbf{P}_h$. The reason is that by setting S_{ij} to -1 , the support of item j will be reduced. Moreover, item j has smaller effect on $\sim\mathbf{P}_h$ than item i .

3. $S_{ij} = 1$, If $\forall \rho \in \mathbf{P}_h$, $\{i, j\} \not\subseteq \rho$ and $\exists \rho' \in \sim\mathbf{P}_h$, such that $\{i, j\} \subseteq \rho'$.
4. $S_{ij} = 0$, otherwise.

Non-Hidden-First Algorithm

Input: $\mathbf{P}_h, \sim\mathbf{P}_h, \mathbf{D}, \mathbf{S}$

Output: \mathbf{D}'

Step 1: Set the values of the entries in \mathbf{S} according to the rules.

Step 2: (matrices multiplication)

For every transaction i in \mathbf{D} do

For $j=1$ to number of items do

If $(D_{ij}=0)$ $D'_{ij}=0$;

Else {

$$\text{Temp} = \sum_{k=1}^{\text{number of items}} D_{ik} * S_{kj}$$

If $(\text{Temp} \geq 1)$ $D'_{ij}=1$;

Else $D'_{ij}=0$; }

However, the sanitized database produced by NHF algorithm may contain sensitive patterns. That is, not all the sensitive patterns can be hidden successfully by applying NHF algorithm.

3.3 HPCME Algorithm

The main idea of HPCME algorithm (Hiding sensitive Patterns Completely with Minimum side Effect on non-sensitive patterns) is to combine the advantages in HF and NHF. All sensitive patterns will be hidden with minimal side effect on non-sensitive patterns.

Restoration probability

Based on HF algorithm, NHF algorithm set proper entries in S to 1 to avoid canceling non-sensitive patterns accidentally. However, some sensitive patterns may be hidden unsuccessfully. Therefore, a new factor restoration probability ($0 \leq p_r \leq 1$) is introduced to decide whether the value of D'_{ij} would follow the multiplication result when the multiplication result is 1 and there exist a $S_{kj} = -1$ ($1 \leq k \leq \text{number of items}$).

A higher value of p_r will let HPCME algorithm tend to preserve non-sensitive patterns, and vice versa. Because our goal is to hide all the sensitive patterns with minimum side effect on non-sensitive patterns, p_r is set to a small value in HCPME algorithm. Moreover, the entries in S are set according to the rules defined in NHF algorithm.

HPCME Algorithm

Input: $P_h, \sim P_h, D, S, p_r$

Output: D'

Step 1: Set the values of the entries in S according to the rules.

Step 2: (matrices multiplication)

For every transaction i in D do

For $j=1$ to number of items do

If ($D_{ij}=0$) $D'_{ij}=0$;

Else {

$$\text{Temp} = \sum_{k=1}^{\text{number of items}} D_{ik} * S_{kj}$$

If ($\text{Temp} \leq 0$) $D'_{ij}=0$

Else {

if ($\exists S_{kj} = -1, 1 \leq k \leq \text{number of items}$)

{ $D'_{ij}=1$ with probability p_r

$D'_{ij}=0$ with probability $1-p_r$ }

else $D'_{ij}=1$ }

4. Performance Evaluation

4.1 The Metrics for Quantifying Performance

Two metrics are introduced to evaluate the effectiveness of our algorithms.

Error 1 : some sensitive patterns can still be discovered after sanitization process. The hiding accuracy is measured by

$$\text{Accuracy} = \frac{\text{number of patterns in } P_h \text{ which are hidden successfully}}{\text{number of patterns in } P_h}$$

A sensitive pattern p_s is said to be hidden successfully if there is not exist a pattern p , such that p can be discovered from D' where p is a subpattern of p_s and p is not a subpattern of any non-sensitive pattern.

Error 2 : some non-sensitive patterns are hidden after sanitization process.

$$\text{Wrongness} = \frac{\text{number of patterns in } \sim P_h \text{ which are disappeared after the sanitization process}}{\text{number of patterns in } \sim P_h}$$

Moreover, overlap rate is defined as follows for evaluating our approach.

$$\text{overlap rate} = \frac{|\text{item}(P_h) \cap \text{item}(\sim P_h)|}{|\text{item}(P_h \cup \sim P_h)|}$$

Where $\text{item}(P)$ denotes the set of items contained by P and $|X|$ denotes the cardinality of set X .

4.2 Experiment Results

The test dataset is generated by the IBM synthetic data generator. The test dataset contains 200 different items, with 100K transactions. Moreover, p_r is set to 0.35 in our experiments.

Figure 2 shows the accuracy of algorithms HF, HPCME and NHF. As shown in the result, HF and

HPCME algorithm approach at 100% accuracy no matter what the values of overlap rate. NHF works like HF and HCME when the overlap rate is low. However, as overlap rate increases, its accuracy decreases.

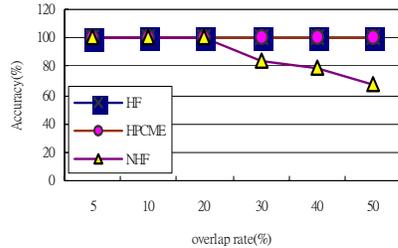


Figure 2. Effect of overlapped rate on Accuracy

Figure 3 shows the wrongness of algorithms for HF, HPCME and NHF. The more the sensitive patterns need to be hidden, the more the entries in CM are set to -1. As a result, non-sensitive patterns are missed easily.

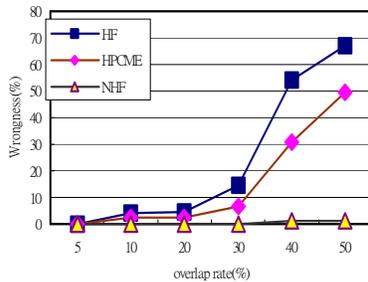


Figure 3. Effect of overlapped rate on Wrongness

5. Conclusion and Future Works

In this paper, a new framework is presented for enhancing privacy in mining frequent patterns. By setting the entries in the sanitization matrix to appropriate values, and multiplying original DB and sanitization matrix, a sanitized database is gotten. According to different settings in sanitization matrix, we bring up three sanitization algorithms for hiding

sensitive patterns successfully or for no legitimate pattern missing.

Acknowledge

This work was partially supported by the National Science Council in Republic of China under the Contract No. 922213E259006

Reference

- [1] M. J. Atallah, E. Bertino, A. Elmagarmid, M. Ibrahim and V. S. Verykios. "Disclosure Limitation of Sensitive Rules". Proceeding of IEEE Knowledge and Data Engineering Exchange Workshop, November 1999.
- [2] R. Agrawal and R. Srikant. Fast algorithms for mining association rules. In *Proc. 1994 Int. Conf. Very Large Data Bases*, 1994.
- [3] R. Agrawal and R. Srikant "Privacy Preserving Data Mining", Proceeding of ACM SIGMOD, 2000.
- [4] R. Agrawal, R. Srikant, A. Evfimievski and J. Gehrke "Privacy Preserving of Association Rules", Proceeding of 8th ACM SIGMOD Conference on Knowledge Discovery and Data Mining (KDD), July 2002.
- [5] A. Evfimievski, J. Gehrke and R. Srikant "Limiting Privacy Breaches in Privacy Preserving Data Mining", Proceeding of ACM PODS, 2003.
- [6] Shariq J. Rizvi and Jayant R. Haritsa, "Maintaining Data Privacy in Association Rule Mining", Proceedings of the 28th VLDB Conference, 2002.
- [7] Stanley R. M. Oliveira, Osmar R. Zaiane "Privacy Preserving Frequent Itemset Mining", IEEE International Conference on Data Mining Workshop on Privacy, Security, and Data Mining, 2002.
- [8] Y. Saygin, V. Verykios and C. Clifton, "Using Unknown to Prevent Discovery of Association Rules", ACM SIGMOD Records, Vol.30, no.4, 2001