

# New framework in Sensitive Rule Hiding

A.S. NAVEENKUMAR

Department of Finance and Computer Applications,  
S.N.R. Sons College,  
Coimbatore .

Dr. M. PUNITHAVALLI

Director,  
Department of Computer Applications,  
Sri Ramakrishna Engineering College Cbe

## Abstract:

Data mining is the process of extracting hidden patterns from data. As more data is gathered, with the amount of data doubling every three years, data mining is becoming an increasingly important tool to transform this data into information. Privacy preserving data mining is a novel research direction in data mining and statistical databases, which has recently been proposed in response to the concerns of preserving personal or sensible information derived from data mining algorithms. There have been two types of privacy proposed concerning data mining. The first type of privacy, called output privacy, is that the data is altered so that the mining result will preserve certain privacy. The second type of privacy, called input privacy, is that the data is manipulated so that the mining result is not affected or minimally affected. For output privacy in hiding association rules, current approaches require hidden rules or patterns to be given in advance. However, to specify hidden rules, entire data mining process needs to be executed. For some applications, only certain sensitive rules that contain sensitive items are required to hide. In this work, an algorithm ISSRH (Increase Support Sensitive Rule Hiding) is proposed, to hide the sensitive rules that contain sensitive items, so that sensitive rules containing specified sensitive items on the right hand side of the rule cannot be inferred through association rule mining.

**Keywords:** Data Mining, Privacy Preserving, Association Rules, Sensitive Rules, Clustering, Minimum Support, Minimum confidence

## 1. Introduction

The concept of Privacy-Preserving has recently been proposed in response to the concerns of preserving personal or sensible information derived from data mining algorithms. Successful applications of data mining have been demonstrated in marketing, business, medical analysis, product control, engineering design, bioinformatics and scientific exploration, among others. The current status in data mining research reveals that one of the current technical challenges is the development of techniques that incorporate security and privacy issues.

Providing security to sensitive data against unauthorized access has been a long term goal for the database security research community and for the government statistical agencies. Recent advances in data mining technologies have increased the disclosure risks of sensitive data.

Hence, the security issue has become, recently, a much more important area of research.

Therefore, in recent years, privacy-preserving data mining has been studied extensively. A number of algorithmic techniques have been designed for privacy-preserving data mining. Most methods use some form of transformation on the data in order to perform the privacy preservation. Typically, such methods reduce the granularity of representation in order to reduce the privacy. This reduction in granularity results in some loss of effectiveness of data management or mining algorithms. This is the natural trade-off between information loss and privacy. Some examples of such techniques are as follows:

- *The randomization method:* The randomization method is a technique for privacy preserving data mining in which noise is added to the data in order to mask the attribute values of records. The noise added is sufficiently large so that individual record values cannot be recovered.

- *The k-anonymity model and l-diversity:* The k-anonymity model was developed because of the possibility of indirect identification of records from public databases. In the k anonymity method, the granularity of data representation is reduced with the use of techniques such as generalization and suppression.

In this work, the sensitive rules are given and the algorithm ISSRH is proposed to modify data in database so that sensitive rules containing specified sensitive items on the right hand side of rule cannot be inferred through association rule mining. The proposed algorithm is based on modifying or perturbing the database transactions so that the confidence of the association rules can be reduced.

## 2. Related Work

In output privacy, given specific rules or patterns to be hidden, many data altering techniques for hiding association, classification and clustering rules have been

proposed. For association rules hiding, two basic approaches have been proposed. The first approach (Saygin , Verykios , Clifton, 2001, Verykios, Elmagarmid , Bertino , Saygin , Dasseni ,2004) hides one rule at a time. It first selects transactions that contain the items in a give rule. It then tries to modify transaction by transaction until the confidence or support of the rule fall below minimum confidence or minimum support. Either removing items from the transaction or inserting new items to the transactions does the modification of transaction. The second approach (Oliveira, Zaiane, 2002 a, Oliveira, Zaiane, 2002 b, Oliveira, Zaiane, 2003 a, Oliveira, Zaiane, 2003 b) deals with groups of restricted patterns or association rules at a time. It first selects the transactions that contain the intersecting patterns of a group of restricted patterns. Depending on the disclosure threshold given by users, it sanitizes a percentage of the selected transactions in order to hide the restricted patterns.

However, both the above approaches require hidden rules or patterns been given in advance.

The work presented here differs from the related work in some aspects are as follows: First, database indexing is performed. Second, correlations among the sensitive rules are considered. Third, avoids the modification in transactions unnecessarily, if the confidence of the sensitive rule gets reduced. Fourth, alter the transactions in the cluster and finally changes can be updated in the database, which reduces the time period of database updating.

### 3. Problem Statement

The problem of mining association rules was introduced in (Agrawal, Imielinski, Swami, 1993). Let  $I = \{i_1, i_2, \dots, i_m\}$  be a set of literals, called items. Given a set of transactions  $D$ , where each transaction  $T$  in  $D$  is a set of items such that  $T \subseteq I$ , an association rule is an expression  $X \Rightarrow Y$  where  $X \subseteq I$ ,  $Y \subseteq I$ , and  $X \cup Y = \pi$ . The  $X$  and  $Y$  are called respectively the body (left hand side) and head (right hand side) of the rule. An example of such a rule is that 90% of customers buy hamburgers also buy Coke. The 90% here is called the confidence of the rule, which means that 90% of transaction that contains  $X$  (hamburgers) also contains  $Y$  (Coke). The confidence is calculated as  $|X \cup Y| / |X|$ , where  $|X|$  is the number of transactions containing  $X$  and  $|X \cup Y|$  is the number of transactions containing both  $X$

and  $Y$ . The notation  $U$  here is not the set union operator. The support of the rule is the percentage of transactions that contain both  $X$  and  $Y$ , which is calculated as  $|X \cup Y| / N$ , where  $N$  is the number of transactions in  $D$ . In other words, the confidence of a rule measures the degree of the correlation between item sets, while the support of a rule measures the significance of the item sets. A typical association rule-mining algorithm first finds all the sets of items that appear frequently enough to be considered significant and then it derives from them the association rules that are strong enough to be considered interesting. The problem of mining association rules is to find all rules that are greater than the user-specified minimum support and minimum confidence.

The objective of data mining is to extract hidden or potentially unknown but interesting rules or patterns from databases.

However, the objective of privacy preserving data mining is to hide certain sensitive information so that they cannot be discovered through data mining techniques (Agrawal, Imielinski, Swami, 1993, Evfimievski, Gehrke , Srikant, 2003).

In this work, an algorithm ISSRH (Increase Support Sensitive Rule Hiding) is proposed, to hide the sensitive rules that contain sensitive items, so that sensitive rules containing specified sensitive items on the right hand side of rule cannot be inferred through association rule mining. More specifically, given a transaction database  $D$ , a minimum support, a minimum confidence and a set of sensitive items  $Y$ , the objective is to minimally modify the database  $D$  such that no sensitive rules containing sensitive items  $Y$  on the right hand side of the rule will be discovered.

### 4. Framework of the approach

Figure 1 shows the framework of the approach that consists of six processes. Initially, indexing is performed in the database. Then association rules are mined from the database. Sensitive items are identified to find the sensitive rules.

Then the sensitive rules are generated. Clustering is performed on the sensitive rules to group the similar items. The rule hiding process is performed and the transactions are updated in the transaction table and finally it is updated in the original database. The main challenge of rule hiding is how to select the items and

transactions to modify. The proposed framework hides the sensitive rules.

## 5. Proposed Algorithm

In order to hide an sensitive rule,  $X \Rightarrow Y$ , it can be either decrease its supports,  $(|X|/N$  or  $|X \cup Y|/N$ ), to be smaller than pre-specified minimum support or its confidence  $(|X \cup Y|/|X|)$  to be smaller than pre-specified minimum confidence.

In the transactions that do not contain both X and Y, to increase the support of X only, the left hand side of the rule, it would reduce the confidence of the rule. In order to hide sensitive rules, when considering hiding sensitive rules with 2 items,  $x \Rightarrow z$ , where z is a sensitive item and x is a single large one item. In theory, association rules may have more specific rules that contain more items, e.g.,  $xY \Rightarrow z$ , where Y is a large item set. However, for such rule to exist, its confidence must be greater than the confidence of  $x \Rightarrow z$ , i.e.,  $\text{conf}(xY \Rightarrow z) > \text{conf}(x \Rightarrow z)$  or  $|xY \cup z| > \text{conf}(x \Rightarrow z) \cdot |xY|$ . For higher confidence rules, such as  $\text{conf}(x \Rightarrow z) = 1$ , there will be no more specific rules. In addition, once the more general rule is hidden, the more specific rule might be hidden as well.

The algorithm tries to decrease the support of the right hand side of the rule.

### Algorithm ISSRH

Input:

1. source database D
2. min support
3. min conf
4. sensitive items Y

Output:

a transformed database D', where rules containing Y on RHS will be hidden.

Step 1: Indexing the transactional database.

Step 2: Generate the association rules.

Step 3: Selecting the Sensitive rules with single antecedent and consequent with the sensitive item in the consequent.(x->y)

Step 4: Constraint based clustering - Clustering the rules with the right hand side has the common item and indexing the rules.

Step 5: Check all the rules in the cluster.

for cluster i= 1 to n

{

do

{

1 select the rule with high confidence

2 find K= mincount.

3 |Tr| = no of sensitive rules that are correlated Such that  $k < |Tr|$

4 Tk = find first k transactions that do not contains the items in the correlated rules

5 sort Tk in ascending order by the number of items.

6 For i= 1 to k

{

7 Choose the first transaction t from Tk.

8 Modify t to support x.  $\text{count}(x)=\text{count}(x) +1$ .

9 Find the confidence of every rule

10 If confidence of all the rules  $< \text{min conf}$

{

11 Update the transactions in the clusters

12 Exit

13 Else

14 Update the transactions in the clusters

}

}

15 unchecked the rules  $< \text{minconf}$

16 }while(rules in cluster[i]==checked)

17 }

Step 4: update database D as transformed database D'.

The algorithm tries to generate the association rule using Agrawal, Imielinski, Swami, 1993). Then it selects the sensitive rules with the sensitive items in the right hand side. Cluster the rules with the common item in the right hand side of the rule (Han, Kamber,2001) and index the rule using (Oliveira, Zaiane, 2003 a).

The mincount= [count(xUy/min conf] – count(x) +1.

(Yi-Hung Wu, Chia-Ming Chiang, and Arbee L.P. Chen, 2007)

Therefore it would take maximum of k no of executions to hide the rule. The rules in every cluster will be hidden.

#### 6. Example

This section shows the example to demonstrate the proposed algorithm to hide the sensitive rules.

The items in database can be represented as a bit vector 1 and 0.( Saygin Y, Verykios V, Clifton C, 2001)

TID Items Items

T1 ABC 111

T2 ABC 111

T3 ABC 111

T4 AB 110

T5 A 100

T6 AC 101

Frequent item sets are generated with minimum support 0.33. Association rules with minimum confidence 0.70 are generated. The rules and the corresponding confidence values are as follows:

C->B0.75, B->C0.75, C->A1.0, B->A1.0, BC->A1.0, AC->B0.75, AB->C0.75, C->AB0.75, and B-> AC0.75.

The item C is considered as sensitive item. The sensitive rule with single antecedent and consequent is , B->C.

The rule is clustered. In the fifth transaction the item B will be added and placed as 1.

Then the database will be updated as

TID Items Items

T1 ABC 011

T2 ABC 011

T3 ABC 111

T4 AB 110

T5 AB 110

T6 AC 101

After updating the database the rules B->C will have the confidence as 0.6, which is less than minimum confidence and hidden. While hiding the rule the rules AB->C, B->AC are also hidden and the rule A->B is generated as side effect.

#### 7. Analysis

This section analyses some of the characteristics of the proposed algorithm. The first characteristic is the database indexing. The indexing helps in reducing the number of scanning of the database. The second characteristic is the time effect. The time taken to scan the database to search the sensitive rules in the database is reduced because of clustering the sensitive rules. The third characteristic is the database effect. The minimum numbers of transactions are modified because of correlation among the sensitive rules. The fourth characteristic is the efficiency of the algorithm. The database is updated after all the rules are hidden that saves the updating time. The fifth characteristic is the transaction effect. The alteration in the transactions are stopped when the confidence of the sensitive rules are reduced than the minimum confidence.

#### 8. Conclusion

In this work, the database privacy problems caused by data mining technology are discussed and the algorithm for hiding sensitive rules is presented. The proposed algorithm here can automatically hide sensitive rule sets. The previous works does not consider the characteristics that are discussed here. Example illustrating the proposed algorithm is given and the characteristics of the algorithm are analyzed. Further the efficiency of the algorithm will be analyzed and improved by reducing the side effects.

#### References

1. Agrawal R, Imielinski T, Swami A (1993) Mining association rules between sets of items in large databases. In: Proceedings of ACM SIGMOD International Conference on Management of Data, Washington DC
2. Dasseni E, Verykios V, Elmagarmid A, Bertino E (2001) Hiding association rules by using confidence and support. In: Proceedings of 4th Information Hiding Workshop, Pittsburgh, PA, pp 369– 383
3. Evfimievski A (2002) Randomization in privacy preserving data mining. SIGKDD Explorations 4(2), Issue 2:43–48
4. Evfimievski A, Gehrke J, Srikant R (2003) Limiting Privacy Breaches in Privacy Preserving Data Mining. PODS 2003, San Diego, CA
5. Evfimievski A, Srikant R, Agrawal R, Gehrke J (2002) Privacy preserving mining of association rules. In: Proceedings of the 8th ACM SIGKDD Int'l Conference on Knowledge Discovery and Data Mining, Edmonton, Canada
6. Han, J., Kamber, M, ( 2001) Data Mining : Concepts and Techniques., Morgan Kaufmann Publishers.
7. Oliveira S, Zaiane O (2002 a) Privacy preserving frequent item set mining. In: Proceedings of IEEE International Conference on Data Mining, November 2002, pp 43–54.
8. Oliveira S, Zaiane O (2002 b) A framework for enforcing privacy in mining frequent patterns. Technical report, TR02-13, Computer Science Department, University of Alberta, Canada
9. Oliveira S, Zaiane O (2003 a) Algorithms for balancing privacy and knowledge discovery in association rule mining. In: Proceedings of 7th International Database Engineering and Applications Symposium (IDEAS03), Hong Kong
10. Oliveira S, Zaiane O (2003 b) Protecting sensitive knowledge by data sanitization. In: Proceedings of IEEE International Conference on Data Mining
11. Saygin Y, Verykios V, Clifton C (2001) Using unknowns to prevent discovery of association rules. SIGMOD Record 30(4):45– 54
12. Verykios V, Elmagarmid A, Bertino E, Saygin Y, Dasseni E (2004) Association Rules Hiding. IEEE Trans Knowledge and Data Eng 16(4):434–447
13. Yi-Hung Wu, Chia-Ming Chiang, and Arbee L.P. Chen, (2007), Hiding Sensitive Association Rules with Limited Side Effects, IEEE Trans Knowledge and Data Eng, Vol. 19, No. 1.

Fig1 : Framework for Approach

