



# Security Using Anonymization and Slicing

Mr.K.Manikanda Prabhu<sup>1</sup>

Assistant Professor, Sri Ramanathan Engineering College, Tamilnadu, India<sup>1</sup>

**ABSTRACT:** Several anonymization techniques, such as generalization and bucketization, have been designed for privacy preserving micro data publishing. Recent work has shown that generalization loses considerable amount of information, especially for high dimensional data. Bucketization, on the other hand, does not prevent membership disclosure and does not apply for data that do not have a clear separation between quasi-identifying attributes and sensitive attributes. In this paper, we present a novel technique called slicing, which partitions the data both horizontally and vertically. We show that slicing preserves better data utility than generalization and can be used for membership disclosure protection. Another important advantage of slicing is that it can handle high-dimensional data. We show how slicing can be used for attribute disclosure protection and develop an efficient algorithm for computing the sliced data that obey the  $l$ -diversity requirement. Our workload experiments confirm that slicing preserves better utility than generalization and is more effective than bucketization in workloads involving the sensitive attribute. Our experiments also demonstrate that slicing can be used to prevent membership disclosure.

**KEYWORDS:** Anonymization techniques, Generalization, Bucketization, Slicing.

## I. INTRODUCTION

Privacy-preserving publishing of microdata has been studied extensively in recent years. Microdata contains records each of which contains information about an individual entity, such as a person, a household, or an organization. Several microdata anonymization techniques have been proposed. The most popular ones are generalization for  $k$ -anonymity and bucketization for  $l$ -diversity. In both approaches, attributes are partitioned into three categories:(1)some attributes are identifiers that can uniquely identify an individual, such as Name or Social Security Number;(2)some attributes are Quasi-Identifiers (QI), which the adversary may already know (possibly from other publicly-available databases) and which, when taken together, can potentially identify an individual, e.g., Birth date, Sex, and Zipcode;(3)some attributes are Sensitive Attributes (SAs), which are unknown to the adversary and are considered sensitive, such as Disease and Salary. There are various situations in which a person might choose to withhold their identity. Acts of charity have been performed anonymously when benefactors do not wish to be acknowledged. A person who feels threatened might attempt to mitigate that threat through anonymity. In certain situations, it is illegal to remain anonymous. In the United States, 24 states have "Stop and identify" statutes that requires persons detained to self-identify when requested by a law enforcement officer. In recent years, due to increase in ability to store personal data about users and the increasing sophistication of data mining algorithms to leverage this information the problem of privacy preserving data mining has become more important. A number of anonymization techniques have been researched in order to perform privacy-preserving data mining. Data anonymization technique for privacy-preserving data publishing has received a lot of attention in recent years. Detailed data (also called as microdata) contains information about a person, a household or an organization. Most popular anonymization techniques are Generalization and Bucketization. There are number of attributes in each record which can be categorized as 1) Identifiers such as Name or Social Security Number are the attributes that can be uniquely identify the individuals. 2) some attributes may be Sensitive Attributes(SAs) such as disease and salary and 3) some may be Quasi-Identifiers (QI) such as zipcode, age, and sex whose values, when taken together, can potentially identify an individual.



## **II. RELATED WORK**

Anonymity is the condition of having one's name or identity unknown or concealed. It serves valuable social purposes and empowers individuals as against institutions by limiting surveillance, but it is also used by wrong doers to hide their actions or avoid accountability the ability to allow anonymous access to services, which avoid tracking of user's personal information and user behaviour such as user location, frequency of a service usage, and so on. If someone sends a file, there may be information on the file that leaves a trail to the sender. The sender's information may be traced from the data logged after the file is sent.

### **A. Anonymity vs. Security**

Anonymity is a very powerful technique for protecting privacy. The decentralized and stateless design of the Internet is particularly suitable for anonymous behavior. Although anonymous actions can ensure privacy, they should not be used as the sole means for ensuring privacy as they also allow for harmful activities, such as spamming, slander, and harmful attacks without fear of reprisal. Security dictates that one should be able to detect and catch individuals conducting illegal behavior, such as hacking, conspiring for terrorist acts, and conducting fraud. Legitimate needs for privacy should be allowed, but the ability to conduct harmful anonymous behavior without responsibility and repercussions in the name of privacy should not.

### **B. Anonymity vs. Privacy**

Privacy and anonymity are not the same. The distinction between privacy and anonymity is clearly seen in an information technology context. Privacy corresponds to being able to send an encrypted e-mail to another recipient. Anonymity corresponds to being able to send the contents of the e-mail in plain, easily readable form but without any information that enables a reader of the message to identify the person who wrote it. Privacy is important when the contents of a message are at issue, whereas anonymity is important when the identity of the author of a message is at issue.

## **III. PROPOSED SYSTEM**

A novel technique called slicing, which partitions the data both horizontally and vertically. We show that slicing preserves better data utility than generalization and can be used for membership disclosure protection. Another important advantage of slicing is that it can handle high-dimensional data. We show how slicing can be used for attribute disclosure protection and develop an efficient algorithm for computing the sliced data that obey the  $\ell$ -diversity requirement. Our workload experiments confirm that slicing preserves better utility than generalization and is more effective than bucketization in workloads involving the sensitive attribute.

Few procedures used to achieve the goal of proposed protocol are follows.

1. A novel data anonymization technique called slicing to improve the current state of the art.
2. Slicing can be effectively used for preventing attribute disclosure, based on the privacy requirement of  $\ell$ -diversity.
3. An efficient algorithm for computing the sliced table that satisfies  $\ell$ -diversity. This algorithm partitions attributes into columns, applies column generalization, and partitions tuples into buckets. Attributes that are highly-correlated are in the same column.
4. The results confirm that slicing preserves much better data utility than generalization. In workloads involving the sensitive attribute, slicing is also more effective than bucketization. In some classification experiments, slicing shows better performance than using the original data (which may overfit the model). These experiments also show the limitations of bucketization in membership disclosure protection and slicing remedies these limitations.

## **IV. SLICING ALGORITHMS**



Slicing first partitions attributes into columns. Each column contains a subset of attributes. This vertically partitions the table. Slicing also partition tuples into buckets. Each bucket contains a subset of tuples.

#### 4.1 Attribute Partition and Columns

An attribute partition consists of several subsets of A, such that each attribute belongs to exactly one subset. Each subset of attributes is called a column. Specifically, let there be columns C1, C2, . . . , Cc, then  $\cup_{i=1}^c C_i = A$  and for any  $1 \leq i_1 \neq i_2 \leq c$ ,  $C_{i_1} \cap C_{i_2} = \emptyset$ . For simplicity of discussion, consider only one sensitive attribute S. If the data contain multiple sensitive attributes, one can either consider them separately or consider their joint distribution. Exactly one of the c columns contains S. Without loss of generality, let the column that contains S be the last column Cc. This column is also called the sensitive column. All other columns {C1, C2, . . . , Cc-1} contain only QI attributes. Our algorithm partitions attributes so that highly correlated attributes are in the same column. This is good for both utility and privacy. In terms of data utility, grouping highly correlated attributes preserves the correlations among those attributes. In terms of privacy, the association of uncorrelated attributes presents higher identification risks than the association of highly correlated attributes because the association of uncorrelated attributes values is much less frequent and thus more identifiable. Therefore, it is better to break the associations between uncorrelated attributes, in order to protect privacy. In this phase, first compute the correlations between pairs of attributes and then cluster attributes based on their correlations.

##### 4.1.1 Measures Of Correlation

Two widely used measures of association are Pearson correlation coefficient and mean square contingency coefficient. Pearson correlation coefficient is used for measuring correlations between two continuous attributes while mean-square contingency coefficient is a chi-square measure of correlation between two categorical attributes. Choose to use the mean-square contingency coefficient because most of our attributes are categorical. Given two attributes A1 and A2 with domains  $\{v1_1; v1_2; \dots; v1_{d1}\}$  and  $\{v2_1; v2_2; \dots; v2_{d2}\}$ , respectively. Their domain sizes are thus d1 and d2, respectively. The mean-square contingency coefficient between A1 and A2 is defined as

$$\phi^2(A_1, A_2) = \frac{1}{\min\{d_1, d_2\} - 1} \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \frac{(f_{ij} - f_i \cdot f_j)^2}{f_i \cdot f_j}$$

Here,  $f_i$  and  $f_j$  are the fraction of occurrences of  $v1_i$  and  $v2_j$  in the data, respectively.  $f_i$  and  $f_j$  is the fraction of occurrences of  $v1_i$  and  $v2_j$  in the data.

#### 4.2 COLUMN GENERALIZATION

Tuples are generalized to satisfy some minimal frequency requirement. We want to point out that column generalization is not an indispensable phase in our algorithm. As shown by Xiao and Tao, bucketization provides the same level of privacy protection as generalization, with respect to attribute disclosure. Although column generalization is not a required phase, it can be useful in several aspects. First, column generalization may be required for identity/membership disclosure protection. If a column value is unique in a column (i.e., the column value appears only once in the column), a tuple with this unique column value can only have one matching bucket. This is not good for privacy protection, as in the case of generalization/bucketization where each tuple can belong to only one equivalence-class/bucket. The main problem is that this unique column value can be identifying. In this case, it would be useful to apply column generalization to ensure that each column value appears with at least some frequency. Second, when column generalization is applied, to achieve the same level of privacy against attribute disclosure, bucket sizes can be smaller. While column generalization may result in information loss, smaller bucket-sizes allow better data utility. Therefore, there is a trade-off between column generalization and tuple partitioning. The trade-off between column generalization and tuple partitioning is the subject of



future work. The algorithms can be applied on the subtable containing only attributes in one column to ensure the anonymity requirement.

### 4.3 TUPLE PARTITIONING

In the tuple partitioning phase, tuples are partitioned into buckets, no generalization is applied to the tuples. Section 4.3.2 gives the description of the tuple-partition algorithm. The algorithm maintains two data structures:

- 1) a queue of buckets Q
- 2) a set of sliced buckets SB.

Initially, Q contains only one bucket which includes all tuples and SB is empty. For each iteration, the algorithm removes a bucket from Q and splits the bucket into two buckets. If the sliced table after the split satisfies  $l$ -diversity, then the algorithm puts the two buckets at the end of the queue Q. Otherwise, we cannot split the bucket anymore and the algorithm puts the bucket into SB (line 7). When Q becomes empty, we have computed the sliced table. The set of sliced buckets is SB (line 8). The main part of the tuple-partition algorithm is to check whether a sliced table satisfies ' $l$ -diversity (line 5).

#### 4.3.1 Algorithm Tuple- Partition

1.  $Q = \{T\}$ ;  $SB = \Phi$ .
2. While Q is not empty
3. remove the first bucket B from Q;  $Q = Q - \{B\}$ .
4. split B into two buckets B1 and B2, as in Mondrian.
5. if  $\text{diversity-check}(T, Q \cup \{B1, B2\} \cup SB, \ell)$
6.  $Q = Q \cup \{B1, B2\}$ .
7. else  $SB = SB \cup \{B\}$ .
8. return SB.

#### 4.3.2 Algorithm Diversity-Check

1. for each tuple  $t \in T$ ,  $L[t] = \Phi$ .
2. for each bucket B in  $T^*$
3. record  $f(v)$  for each column value v in bucket B.
4. for each tuple  $t \in T$
5. calculate  $p(t, B)$  and find  $D(t, B)$ .
6.  $L[t] = L[t] \cup \{(p(t, B), D(t, B))\}$ .
7. for each tuple  $t \in T$
8. Calculate  $p(t, s)$  for each s based on  $L[t]$ .
9. if  $p(t, s) \geq 1/\ell$ , return false.
10. return true.

## V. PROTOCOL OVERVIEW

The modules used in this paper are as follows

Original Data, Generalized Data, Bucketized Data. Multiset-based, Generalization Data, One-attribute-per-Column Slicing Data, Sliced Data.

### 5.1 Original Data



## International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol.2, Special Issue 1, March 2014

### Proceedings of International Conference On Global Innovations In Computing Technology (ICGICT'14)

Organized by

Department of CSE, JayShriram Group Of Institutions, Tirupur, Tamilnadu, India on 6<sup>th</sup> & 7<sup>th</sup> March 2014

We conduct extensive workload experiments. Our results confirm that slicing preserves much better data utility than generalization. In workloads involving the sensitive attribute, slicing is also more effective than bucketization. In some classification experiments, slicing shows better performance than using the original data.

#### 5.2 Generalized Data

Generalized Data, in order to perform data analysis or data mining tasks on the generalized table, the data analyst has to make the uniform distribution assumption that every value in a generalized interval/set is equally possible, as no other distribution assumption can be justified. This significantly reduces the data utility of the generalized data.

Generalization transforms the QI-values in each bucket into “less specific but semantically consistent” values so that tuples in the same bucket cannot be distinguished by their QI values.

#### 5.2 Bucketized Data

We show the effectiveness of slicing in membership disclosure protection. For this purpose, we count the number of fake tuples in the sliced data. We also compare the number of matching buckets for original tuples and that for fake tuples. Our experiment results show that bucketization does not prevent membership disclosure as almost every tuple is uniquely identifiable in the bucketized data.

#### 5.3 Multiset-Based Generalization Data

We observe that this multiset-based generalization is equivalent to a trivial slicing scheme where each column contains exactly one attribute, because both approaches preserve the exact values in each attribute but break the association between them within one bucket.

Because each Generalized separately, correlations between different attributes are lost. In order to study attribute correlations on the generalized table, the data analyst has to assume that every possible combination of attribute values is equally possible. This is an inherent problem of generalization that prevents effective analysis of attribute correlations.

#### 5.4 One-Attribute-Per-Column Slicing Data

We observe that while one-attribute-per-column slicing preserves attribute distributional information, it does not preserve attribute correlation, because each attribute is in its own column. In slicing, one group correlated attributes together in one column and preserves their correlation. For example, in the sliced table shown in Table correlations between Age and Sex and correlations between Zip code and Disease are preserved. In fact, the sliced table encodes the same amount of information as the original data with regard to correlations between attributes in the same column.

#### 5.5 Sliced Data

Another important advantage of slicing is its ability to handle high-dimensional data. By partitioning attributes into columns, slicing reduces the dimensionality of the data. Each column of the table can be viewed as a sub-table with a lower dimensionality. Slicing is also different from the approach of publishing multiple independent sub-tables in that these sub-tables are linked by the buckets in slicing. We now present an efficient slicing algorithm to achieve  $\ell$ -diverse slicing. Given a microdata table  $T$  and two parameters  $c$  and  $\ell$ , the algorithm computes the sliced table that consists of  $c$  columns and satisfies the privacy requirement of  $\ell$ -diversity.

## VI. CONCLUSION



This paper presents a new approach called slicing to privacy-preserving microdata publishing. Slicing overcomes the limitations of generalization and bucketization and preserves better utility while protecting against privacy threats. We illustrate how to use slicing to prevent attribute disclosure and membership disclosure. Our experiments show that slicing preserves better data utility than generalization and is more effective than bucketization in workloads involving the sensitive attribute. The general methodology proposed by this work is that: before anonymizing the data, one can analyze the data characteristics and use these characteristics in data anonymization. The rationale is that one can design better data anonymization techniques when we know the data better. We show that attribute correlations can be used for privacy attacks. This work motivates several directions for future research. First, in this paper, we consider slicing where each attribute is in exactly one column. An extension is the notion of over-lapping slicing, which duplicates an attribute in more than one column. These releases more attribute correlations. For example, in Table sliced data, one could choose to include the Disease attribute also in the first column. That is, the two columns are {Age,Sex,Disease} and {Zipcode,Disease}. This could provide better data utility, but the privacy implications need to be carefully studied and understood. It is interesting to study the trade off between privacy and utility. Second, we plan to study membership disclosure protection in more details. Our experiments show that random grouping is not very effective. Third, slicing is a promising technique for handling high dimensional data. By partitioning attributes into columns, we protect privacy by breaking the association of uncorrelated attributes and preserve data utility by preserving the association between highly-correlated attributes. For example, slicing can be used for anonymizing transaction databases, which has been studied recently in. Finally, while a number of anonymization techniques have been designed, it remains an open problem on how to use the anonymized data. In our experiments, we randomly generate the associations between column values of a bucket.

## REFERENCES

- [1] J. Brickell and V. Shmatikov. The cost of privacy: destruction of data-mining utility in anonymized data publishing. In KDD, pages 70–78, 2008.
- [2] B.-C. Chen, R. Ramakrishnan, and K. LeFevre. Privacy skyline: Privacy with multidimensional adversarial knowledge. In VLDB, pages 770–781, 2007.
- [3] K. LeFevre, D. J. DeWitt, and R. Ramakrishnan. Workload-aware anonymization. In KDD, pages 277–286, 2006.
- [4] J. Li, Y. Tao, and X. Xiao. Preservation of proximity privacy in publishing numerical sensitive data. In SIGMOD, pages 473–486, 2008.
- [5] N. Li, T. Li, and S. Venkatasubramanian.  $t$ -closeness: Privacy beyond  $k$ -anonymity and  $\ell$ -diversity. In ICDE, pages 106–115, 2007.
- [6] T. Li and N. Li. Injector: Mining background knowledge for data anonymization. In ICDE, pages 446–455, 2008.
- [7] T. Li and N. Li. On the trade off between privacy and utility in data publishing. In KDD, pages 517–526, 2009.
- [8] H. Cramt'er. Mathematical Methods of Statistics. Princeton, 1948.
- [9] I. Dinur and K. Nissim. Revealing information while preserving privacy. In PODS, pages 202–210, 2003.
- [10] C. Dwork. Differential privacy. In ICALP, pages 1–12, 2006.
- [11] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In TCC, pages 265–284, 2006.
- [12] J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. TOMS, 3(3):209–226, 1977.
- [13] B. C. M. Fung, K. Wang, and P. S. Yu. Top-down specialization for information and privacy preservation. In ICDE, pages 205–216, 2005.
- [14] G. Ghinita, Y. Tao, and P. Kalnis. On the anonymization of sparse high-dimensional data. In ICDE, pages 715–724, 2008.