

PRIVACY PRESERVATION FOR MULTIPLE SENSITIVE ATTRIBUTES

Yang Ye

Institute forTheoretical Computer Science, Tsinghua University, Beijing 100084, China
yey05@mails.tsinghua.edu.cn

Dapeng Lv, Yu Liu, Jianhua Feng

Dept. of Computer Science, Tsinghua University, Beijing 100084, China
{liyuu-05,lvdp05}@mails.tsinghua.edu.cn, fengjh@tsinghua.edu.cn

I .Background:

- Three types of attributes in the mcrodata:**
 - Identifying Attributes:
Name, Social Security Number
 - Quasi-identifying (QI) Attributes:
Date of birth, Zip code and Gender
 - Sensitive Attributes:
Disease, Salary
- Objectives of Privacy Preservation:**
 - Prevent direct or indirect disclosure of sensitive values.
 - Enable the researcher to effectively investigate the relationship between sensitive attributes and other attributes.
- The prevalent privacy preservation technique: anonymization**
 - Eliminate Identifying Attributes.
 - Generalization on QI Attributes to form QI groups.
 - Anonymization principles on QI groups:
k-anonymity, l-diversity, etc.
- Anonymization principles: bound the strength of privacy preservation**
 - K-anonymity(against link attack):
Each QI group with size at least k .
 - L-diversity(against homogenous and background attack):
Each QI group contains at least l “well-represented” sensitive values.

II .Two Different Cases of Privacy Preservation:

- The SSA Case(the work of the state of the art):**
Each tuple contains one single sensitive attribute
- The MSA Case(our Work):**
Each tuple contains multiple sensitive attributes

vi. Decompose in the SSA case:

- The **Diversity Parameter l**: resembles the concept in l-diversity
- The **Group-Forming Method : Largest---l Method**
 - Bucketization :**
Tuples with identical sensitive attributes are placed in the same bucket;
Bi: the i-th largest buckets
 $|B_i| = n_i$
 - Group Forming:**
In each iteration of group forming, one tuple is removed from each of the l largest buckets to form a new SA-group

after one iteration, the size of some buckets will be changed. So in the beginning of every iteration, the buckets are sorted according to their sizes

The result of decompose on **Table II** is depicted in **Table IV**

- Intensive Study on **Largest---l Method**:

Theorem 1: The **Largest-l** group forming method creates as many groups as possible

DEFINITION 3 (l-Property) : The original data distribution satisfies l-Property iff

- $n_i / n \leq 1 / l$
- $n = k \cdot l$

Theorem 2: If the original data distribution satisfies l-Property, then after the **Largest-l** Group Forming method, no tuple will be left.

Corollary 1 : If the original data distribution satisfies l-Property (1) while does not satisfy l-Property (2), then after the **Largest-l** Group Forming procedure , there will be only one tuple in the non-empty buckets.

Corollary 2: The optimal assignment of diversity parameter l is $\lfloor n/n_1 \rfloor$

III. The Running Example

TABLE I THE MICRODATA TABLE						TABLE II PART OF A VOTE REGISTER LIST				TABLE III THE GENERALIZED TABLE					
Tuple#	Sex	Zip	Birth.	Occ.	Sal.	Name	Sex	Zip	Birth.	#	Sex	Zip	Birth.	Occ.	Sal.
1(Alice)	F	10078	1988/04/17	nurse	1	Alice	F	10078	1988/04/17	1	*	1007*	1983-88	nurse	1
2(Betty)	F	10077	1984/03/21	nurse	4	Betty	F	10077	1984/03/21	2	*	1007*	1983-88	nurse	4
3(Carl)	M	10076	1985/03/01	police	8	Carl	M	10076	1985/03/01	3	*	1007*	1983-88	police	8
4(Diana)	F	10075	1983/02/14	cook	9	Diana	F	10075	1983/02/14	4	*	1007*	1983-88	cook	9
5(Ella)	F	10085	1962/10/03	actor	2	Ella	F	10085	1962/10/03	5	*	*008*	1958-88	actor	2
6(Finch)	M	10085	1988/11/04	actor	7	Finch	M	10085	1988/11/04	6	*	*008*	1958-88	actor	7
7(Gavin)	M	20086	1958/06/06	clerk	8	Gavin	M	20086	1958/06/06	7	*	*008*	1958-88	clerk	8
8(Helen)	F	20087	1960/07/11	clerk	2	Helen	F	20087	1960/07/11	8	*	*008*	1958-88	clerk	2

- Here the value i for **salary** means monthly income is **1000i-----1000(i+1)** dollars.
- In this example, “**Sex**,” “**Zip**” and “**Birth.**” are treated as QI attributes while “**Occ.**” and “**Salary**” are sensitive attributes.
- Through **generalization**, a generalized table, TABLE III, is formed, Which contains two QI-groups.
- Each group contains 4 tuples, therefore, TABLE III satisfies **4-anonymity**.
- Besides, for either sensitive attribute “**Occ.**” or “**Salary**”, the first group contains at least 3 different values.
Group 1 satisfies **3-diversity** for “Occ.” and “Salary” respectively.

IV. New Privacy Risks in the MSA cases

- Chain Attack:**
If an adversary locate **Carl** in the first group and he previously knows Carl’ **occupation** is “police”, he will obtain Diana’s “**Salary**” information of “**8000-9000**” with full confidence. This kind of attack is termed “**link attack**”.
- Exclusion Attack:**
If an adversary locates **Carl** in the first group and although he does not know Carl’s occupation, but he can conclude Carl is not a **nurse** because nurse is an job for women. So, the adversary knows Carl corresponds to the 3rd or the 4th tuple, so he can deduce Carl’s “**Salary**” information is “**8000-10000**”. This kind of attack is termed “**exclusion attack**”.
- The Mechanism of new Privacy Risks in Anonymization in the MSA cases.**
In each group, the distribution of “Salary” attribute for each value of “Occupation” l lacks diversity.
- A QI group like group **2** may be satisfactory, in which each occupation value corresponds to 2 different salary values. However, we can prove to obtain such groups to cover the whole table is indeed impractical.

V . Our solution: Decompose

- Decompose in the **SSA** case:
Resembles “**Anatomy**” and “**Permutation**” Table will be partitioned into “**SA-groups**”

DEFINITION 1 (SA-Group)

A SA-grp contains tuples with untransformed QI values and each tuple is associated with the union of these tuples’ original sensitive values.

- Choose one Sensitive attribute as “**Primary Sensitive Attribute**” and form SA-groups with this attribute
The **Group-Forming Method**: To Be Discussed Latter.

DEFINITION 2 (Primary Sensitive Attribute)

In the MSA case, the primary sensitive attribute is the sensitive attribute chosen by the publisher, according to which SA-groups are formed.

- Extending to Multiple Sensitive Attributes: By possibly **adding noise** in other attributes.
The **Noise-Adding Method**: To Be Discussed Latter.

VII. Extending Decompose to the MSA case.

- Form SA-groups according to the **Primary Sensitive Attribute**
- For each other group and each other sensitive attribute, unite up the original values, reduplicated values are counted just once, because multiple counts just increase privacy disclosure risk, as shown in **Group 1** of **Table V** (salary value 8)
- Possibly in combination with **Adding Noise**
 - In **Group 1** of **Table V** , there are only 3 distinct values for “salary”. However, the optimal assignment of diversity parameter for “salary” is $l=8/2=4$.
 - Noise values are not arbitrarily chosen, in fact, **4** and **7** are allowed noise values while **9** is not.
 - Because by linking **Group 1**’s **Occupation** values with the **sensitive table**, the adversary can deduce **9** cannot be a salary value of Group 1. *Detailed implementation of Adding Noise is neglected for lack of space*
 - The final publishing of Decompose is shown in Table VI, where we choose 4 as the noise value

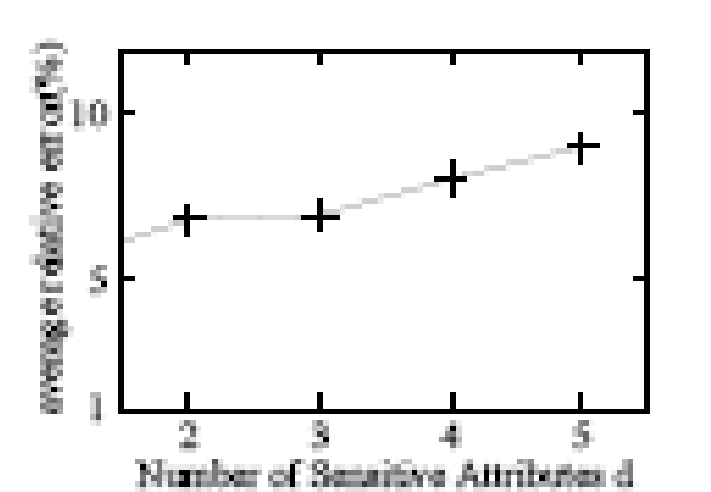
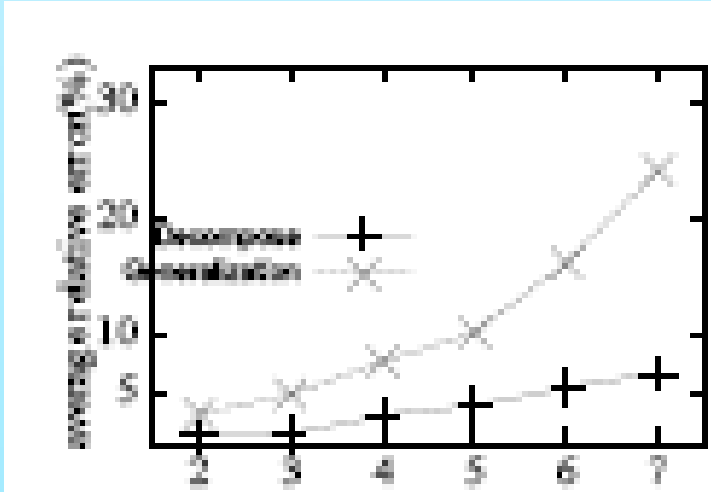
TABLE VI
THE FINAL PUBLISHING OF DECOMPOSE

The Sensitive Table		The Decomposed Table after Adding Noise					
Occ.	Sal.	Grp#	Sex	Zip	Birth.	Occ.	Sal.
nurse	1	1	F	10078	1988/04/17	police	1
nurse	4		F	10085	1962/10/03	nurse	2
police	8		M	20086	1958/06/06	actor	4
cook	9		M	10076	1985/03/01	clerk	8
actor	2	2	F	10077	1984/03/21	nurse	2
actor	7		M	10085	1988/11/04	actor	4
clerk	8		F	10075	1983/02/14	cook	7
clerk	2		F	20087	1960/07/11	clerk	9

VIII. Experiments

- Setting
Real Database : **Adult** (Downloaded at :<http://www.ics.uci.edu/mmllearn/mlrepository.html>)
- The measurement : *average relative error* in answering aggregate query
 $relative\ error = |act - est| / act$
act = actual result derived from the microdata
est = the estimate computed from the published table.
- Aggregate Query:
SELECT COUNT(*) FROM Unknown-Microdata
WHERE pred(A 1) AND ... AND pred(A q) AND pred(S 1) ... AND pred(S d)

pred(A) of the form (A = x1 OR A = x2 OR ... OR A = xb)
- Comparison between Decompose in SSA and Generalization (Parameter l)
- Decompose for multiple sensitive attributes (Parameter d: number of)



IX . Acknowledgement

This work was Supported in part by the National Natural Science Foundation of China Grant 60553001, 60573094, the National Basic Research Program of China Grant 2007CB807900, 2007CB807901, the National High Technology Development 863 Program of China under Grant No.2007AA01Z152 and 2006AA01A101, the National Grand Fundamental Research 973 Program of China under Grant No.2006CB303103, and Basic Research Foundation of Tsinghua National Laboratory for Information Science and Technology (TNList).



清华大学
Tsinghua University