_____

# Slicing: A New Access to Isolation Preserving Data Publishing

Swapnil Y. Gursale
Dept. of Computer Engineering
PVG's College of Engineering
Nasik, India
*syg2590@gmail.com*

Ravindra S. Avhad
Dept. of Computer Engineering
PVG's College of Engineering
Nasik, India
*raviavhad12@gmail.com*

Rajesh K. Pagare
Dept. of Computer Engineering
PVG's College of Engineering
Nasik, India
*rajesh.pagare290@gmail.com*

*Abstract—* There are several anonymizing techniques like Abstraction, Containerization for isolation preserving small data publishing. The Abstraction loses amount of information for high spatial data. Containerization does not avoid enrollment acknowledgment and does not give clear separation between aspects. We are presenting a technique called slicing which partitions the data both horizontally and vertically. We also show that slicing conserves better data service than abstraction and can be used for enrollment acknowledgment conservation. One more important advantage of slicing is that it can handle high-spatial data. Slicing can be used for aspect acknowledgment conservation and establishing an efficient algorithm for computing the sliced data that obey the l-diversity concern. Our experiments confirm that slicing conserves better service than abstraction and is more active than containerization in loads affecting the conscious aspect and also demonstrate that slicing can be used to avoid enrollment acknowledgment.

*Keywords-* *Isolation securities, knowledge publishing, knowledge security*

_____******_____

## I.    INTRODUCTION

### A.    *Problem Definition and Relevant Theory :*

We are studying Isolation preserving of small data from so many years. Small data contains information of a person or organization. There are many techniques proposed.

In that most popular ones are abstraction for [7], [8] k-anonymity[8] and containerization [9], [6], [3] for ℓ-diversity [5]. In both techniques, aspects are divided into three categories: (1) some aspects are differently detectable an individual, such as Name or ID; (2) some aspects are Quasi-Attributes (QI), Quasi-attributes are pieces of information that are not of themselves different attributes, but are ampelly well corrospondance with an entity that they can be combined with other quasi-attributes to create a different attribute. e.g., Birth-date, Sex, and Zipcode; (3) some aspects are Conscious Aspects (SAs), uasi-conscious aspects are not conscious by themselves, but some code or their mixs maybe linked to foreign knowledge to reveal indirect conscious information of an individual. such as Disease and Salary. In both abstraction and containerization, one first removes attributes from the data and then divide records into containers. There is diffrence between two techniques. Abstraction transforms the QI-code in each container into "less specific but semantically dependable" code so that records in the same container cannot be distinguished by their QI code. In containerization, one separates the SAs from the QIs by randomly permuting the SA code in each container. The anonymized data lies of a fixed of containers with permuted conscious aspect code.

### B.    *Motivation of Slicing :*

In abstraction [1], [2], [9] for k anonymity losses large amount of data, specially for high-spatial data. This reasons are: First, abstraction for k-anonymity does not accept high impact of spatiality. So the abstraction to be active, records in the same partitions must be close to each other so that generalizing the records would not lose so much information.

This mainly cuts the data service of the generalized data. Second, because each aspect is generalized disparately, correlations between disparate aspects are lost. To study aspect correlations on the generalized table, the data analyser has to assume that every possible mix of aspect code must be equally possible. This problem of abstraction that avoids active analysis of aspect correlations. While containerization [9], [6], [3] has better data utilization than abstraction but it has several limitations. Because Firstly, containerization does not avoid enrollment acknowledgment[27]. Second, containerization requires a clear separation between QIs and SAs. So in many data fixeds, it is not clear which aspects are QIs and which are SAs. Third, by separation of the conscious aspect from the QI aspects, containerization breaks the aspect correlations between the QIs and the SAs.

In this paper, we announce a new data pseudonyms technique called slicing to improve the retrieval of data. Slicing partitions the data fixed both vertically and horizontally. Vertical dividation is done by combination of aspects into files based on the correlations among the aspects. Each file contains a subfixed of aspects that are highly corrospondance. Horizontal partitioning is done by combination records into containers. Finally, inner each container, code in each file are randomly sorted to break the linking between disparate files.

**2003**

_____

The basic idea behind slicing is to break the association short files, but to preserve the association inner each file. This cuts the spatiality of the data and conserves better service than abstraction and containerization. The service is preserved because it groups highly corrospondance aspects in sync, and conserves the correlations between such aspects. Slicing protects isolation because it breaks the associations between uncorrospondance aspects, which are not frequent and thus identifying.

### C. Contributions & Organization :

Our contributions include the coming. First, we announce slicing as a new technique for isolation preserving data publishing. Slicing has some advantages when compared with abstraction and containerization. It saves better data service than abstraction. It saves more aspect correlations with the SAs than containerization. It also handle high-spatial data and data without a clear separation of QIs and SAs. Second, we show that slicing can be actively used for avoiding aspect acknowledgment, based on the isolation concern of '-diversity. We announce a notion called '- differing slicing, which assure that the attacker cannot learn the conscious value of any individual with a chance greater than 1='. Third, we establishing an efficient algorithm for computing the sliced table that satisfies '-diversity. Our algorithm partitions aspects into files, applies file abstraction, and partitions tuples into containers. Aspects that are highly corrospondance are in the same file; this conserves the correlations between such aspects. The linkings between uncorrospondance aspects are broken; this provides better isolation as the linkings between such aspects are less frequent and potentially identifying.

## II. TYPE STYLE AND FONTS

We now present an efficient slicing algorithm to achieve '- differing slicing. Given a small data table T and two constants c and ', the algorithm computes the sliced table that lies of c files and satisfies the isolation concern of '-diversity. Our algorithm lies of three phases: aspect partitioning,file abstraction, and tuple partitioning. We now explains the three phases.

### A. Aspect Partitioning

Our algorithm partitions aspects so that highly corrospondance aspects are in the same file. This is good for both service and isolation. In terms of data service, combination highly corrospondance aspects conserves the correlations among those aspects. In terms of isolation, the association of uncorre-lated aspects presents higher identification risks than the association of highly corrospondance aspects because the association of uncorrospondance aspect code is much less frequent and thus more detectable. Therefore, it is better to break the associations between uncorrospondance aspects, in order to protect isolation.

*Algorithm tuple partition(T,l):*

1. Q={T};SB=$.
2. While Q is not empty
3. Remove the first container B from Q;Q=Q-{B}.
4. Split B into two containers B1 and B2.,as in Mondrian.
5. If diversity-check(T,Q UNION {B1,B2}UNION SB,l )
6. Q=Q UNION {B1,B2}
7. Else SB=SB UNION {B}
8. Return SB

### B. Coloumn Abstraction

In this second phase, tuples are generalized to satisfy some minimal density concern. We want to mark out that file abstraction is not an indispensable phase in our algorithm. As own by Xiao and Tao[9], containerization provides the same level of isolation conservation as abstraction, with respect to aspect acknowledgment.

### C. Tuple Partitioning

In the records partitioning stage, records are divided into containers. We modify the Mondrian [4] algorithm for tuple partition. Unlike Mondrian k-anonymity, no abstraction is applied to the records; we use Mondrian for the purpose of dividing records into containers.

*Algorithm diversity-check(T,T*,l):*

1. for each tuple t belongs T,L[t]=$
2. for each container B in T*
3. record f(v) for each file value v in container B.
4. for each tuple t belongs T
5. calculate p(t,B) and find D(t,B).
6. L[t]=L[t] UNION {<p(t,B),D(t,B)>}.
7. for each tuple t belongs T
8. calculate p(t,s) for each s based on L[t].
9. if p(t,s)>=1/l.return false
10. return true.

The main part of the tuple-partition algorithm is to check whether a sliced table satisfies '-diversity . For each tuple t, the algorithm maintains a list of statistics L½t_aboutt's matching containers. Each element in the list L½t_ contains statistics about one matching container B: the matching chance pðt;BÞ and the distribution of candidate conscious code Dðt;BÞ.

The algorithm first parts one browse of each container B(lines 2 to 3) to record the density fðvÞ of each file value v in

container B. Then, the algorithm parts one browse of each tuple t in the table T (lines 4 to 6) to find out all tuples that match B and record their matching chance pðt;BÞ and the distribution of candidate conscious code Dðt;BÞ, which are added to the list L½t_ (line 6). At the end of line 6, we accept access, for

each tuple t, the list of statistics L½t_about its matching containers. A final browse of the tuples in Twill compute the pðt; sÞ

code based on the act of total chance.

### D. *Original Data*

We conduct extensive load experiments. Our reactions confirm that slicing conserves much better data service than abstraction. In loads affecting the conscious aspect, slicing is also more active than containerization. In some allocation experiments, slicing shows better achevement than using the original data.

| Account _id | Transaction_no | Phone_no | Balance |
|---|---|---|---|
| 5001 | 49000 | 8109119061 | 10000 |
| 5002 | 49081 | 8810034678 | 20000 |
| 5003 | 59002 | 8743560912 | 5000 |
| 5004 | 39003 | 8754348988 | 15000 |
| 5005 | 79004 | 7398797587 | 30000 |
| 5006 | 67893 | 7654239089 | 22000 |

Table 3.4 Original Data

### E. *Sliced Data*

Another important advantage of slicing is its ability to handle high-spatial data. By partitioning aspects into files, slicing cuts the spatiality of the data. Each file of the table can be considered as a sub-table with a lower spatiality. Slicing is also disparate from the access of publishing multiple independent sub-tables in that these sub-tables are linked by the containers in slicing.

| (Account_id,Transaction_no) | (Phone_no,Balance) |
|---|---|
| (5001,49000) | (8743560912,20000) |
| (5002,49081) | (8810034678,5000) |
| (5003,59002) | (8109119061,10000) |
| (5004,39003) | (7654239089,15000) |
| (5005,7004) | (7398797587,22000) |
| (5006,67893) | (8754348988,30000) |

Table 3.5 Sliced Data

### III.    ENROLLMENT ACKNOWLEDGMENT CONSERVATION

In this   we check how slicing can do enrollment acknowledgment conservation.

### A. *Containerization.*

Let us first check how can attacker can it interfere enrollment data from containerization. Because containerization gives each records pairing of QI code in  original form and most of individuals can be differently identified by the QI code, the attacker can get  the enrollment of an individual in the original record by checking whether the individual's pairing of QI code displays in the releasing data.

### B.Slicing

It gives conservation against enrollment acknowledgment because QI aspects are divided into disparate files and correlations among disparate files inner each container are being  broken. Consider the sliced table in Table 1f. The table has two files. The first container is reactioned from four records; we call them the original records. The container matches alin sync 42¼416 records, 4 original records and 12 that do not appear in the original table. We call these 12 records duplicate records. Given a record, if it has no matching parition in the sliced table, then we ensure that the records is not in the original table. So,even if a record has one or more matching partition, one would not tell whether the record is in the original table, because it can  be a duplicate record. We define two quantitative measures for the degree of enrollment conservation provided by slicing. The first one is the duplicate-original ratio (FOR), which could defined as the number of duplicate records divided by the number of original records. So, the bigger the FOR, the maximum enrollment conservation is given. A sliced container of size k can enhancingly  match kc records, having k original records and kc-k duplicate records; So, the FOR isk c-1-1. When one has chosed a minimum threshold for the FOR, one can part k and c to satisfy the threshold. The second measure is to part the number of pairing partitions for original records and that for duplicate records. If they are same enough, enrollment data is secured because the attacker cannot disparateiate original records from duplictae records. So the main aim of this paper is aspect acknowledgment, we do not imply to give a comprehensive analysis for enrollment acknowledgment conservation.

### C.    ABSTRACTION

By generalizing aspect1 code into "less-specific but semantically dependable code," abstraction provides some conservation against enrollment acknowledgment.The abstraction alone (e.g., used withk-anonymity) may loose enrollment data if the target individual is the only possible pair for a generalized record. The intuition is same to our rationale of duplicate records. If a generalized record does not announce duplicate records (i.e., none of the other mixs of code are reasonable), there will be only one original records that matches with the generalized records and the enrollment information can still be interfered.The major problem is that it

can be difficult to show the backdrop table and in some phases the data publisher may not accept such a backdrop table. And the conservation against enrollment acknowledgment depends on the superior of the backdrop table. Therefore, with careful pseudonyms, abstraction can offer some level of enrollment acknowledgment conservation.

References

[1] C. Aggarwal, "On k-Anonymity and the Curse of Dimension-ality," Proc. Int'l Conf. Very Large Data Bases (VLDB), pp. 901-909, 2005.

[2] D. Kifer and J. Gehrke, "Injecting Service into Anonymized Data Fixeds," ACM SIGMOD Int'l Conf. Management of Data (SIGMOD), pp. 217- 228, 2006.

[3] N. Koudas, D. Srivastava, T. Yu, and Q. Zhang, "Aggregate Query Answering on Anonymized Tables," Proc. IEEE 23rd Int'l Conf. Data Eng. (ICDE), pp. 116-125, 2007.

[4] K.LeFevre, D. DeWitt, and R. Ramakrishnan, "Mondrian Multi-spatial k-Anonymity," Proc. Int'l Conf. Data Eng. (ICDE), 25, 2006.

[5] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkitasubrama-niam, "'-Diversity: Isolation Beyond k-Anonymity," Proc. Int'l Conf. Data Eng. (ICDE), p. 24, 2006.

[6] D.J. Martin, D. Kifer, A. Machanavajjhala, J. Gehrke, and J.Y. Halpern, "Worst-Case Backdrop Knowledge for Isolation-Preserving Data Publishing," Proc. IEEE 23rd Int'l Conf. Data Eng. (ICDE), pp. 126-135, 2007.

[7] P. Samarati, "Protecting Respondent's Isolation in Small data Release," IEEE Trans. Knowledge and Data Eng., vol. 13, no. 6, 1010-1027, Nov./Dec. 2001.

[8] L. Sweeney, "k-Anonymity: A Model for Protecting Isolation," Int'l J. Unsomety Fuzziness and Knowledge-Based Systems, vol. 10, no. 5, 557-570, 2002.