

Privacy Preserving on User Profiles Using TSSA

¹ Ramya Lavanya Parsapogu, ² K.Ramesh Babu

¹Mtech, CHALAPATHI INSTITUTE OF ENGINEERING AND TECHNOLOGY, CHALAPATHI NAGAR, LAM, GUNTUR, AP, INDIA.

²Assistant Professor, CHALAPATHI INSTITUTE OF ENGINEERING AND TECHNOLOGY, CHALAPATHI NAGAR, LAM, GUNTUR, AP, INDIA.

Abstract: Information anonymization is one key part of Micro information divulgements as they empower strategy creators to break down the choice results of issues impacting the business there by affecting the future course of activities. Security is a key issue here on the grounds that improper exposure of certain information stakes will hurt the prospects. Former methodologies of information anonymization, for example, generalization and bucketization (determined by k-obscurity, l-differences) have been intended for protection safeguarding micro information distributed which have a few impediments like Generalization's failure to handle high dimensional information and Bucketization disappointment to keep up clear partition between semi recognizing qualities and touchy characteristics incited the advancement of a novel system called Slicing, which segments the information both evenly and vertically. Albeit Slicing attains better information utility and secrecy contrasted with earlier procedures, its touchy property exposures are focused around arbitrary gathering, which is not extremely viable as haphazardly creating the relationship between section estimations of a basin altogether brings down information utility. In this manner, we propose to supplant irregular gathering with more compelling tuple gathering calculations, for example, Tuple Space Search calculation focused around

hashing systems. The figured and acquired cut information from high dimensional touchy properties focused around the proposed system offers noteworthy execution climb. A possible reasonable usage on dynamic information approves our case.

Index Terms: Privacy Preserving, Data Anonymization, Slicing, Tuple Grouping Method.

I. INTRODUCTION

Information mining that is at times otherwise called Knowledge Discovery Data (KDD) is the procedure of breaking down information from alternate points of view and outlining it into valuable data. Information mining is the concentrating the significant data from the extensive information sets, for example, information stockroom, Micro information holds records each of which holds data about an individual substance. Micro data hold records each of which holds data about an individual element. Numerous micro data anonymization procedures have been proposed and the most famous ones are generalization with k-secrecy and bucketization with l differences. For protection in Micro data distributed a novel method called cutting is utilized that the parcels the information both evenly and vertically. Cutting jam preferable information utility over generalization and might be utilized for participation revelation security. It can

deal with high dimensional information. A finer framework is obliged that can that can with stand high dimensional information taking care of and delicate characteristic divulgence disappointments. These quasi-identifiers are situated of qualities are those that in mix might be joined with the outer data to reidentify. These are three classes of characteristics in microdata. On account of both anonymization strategies, first identifiers are expelled from the information and afterward segment the tuple's into basins.

In generalization, converts the semi recognizing values in each one can into less particular and semantically consistent so that tuple's in the same pail can't be recognized by their QI values. One divides the SA values from the QI values by arbitrarily permuting the SA values in the basin in the bucketization. The anonymized information comprise of a set of pails with permuted touchy property estimations. Existing works for the most part considers datasets with a solitary touchy characteristic while persistent information comprises various delicate traits, for example, determination and treatment.

Information cutting can likewise be utilized to avert participation divulgence and is proficient for high dimensional information and jelly better information utility. We present a novel information anonymization procedure called cutting to enhance the current state of the craft. Information has been apportioned evenly and vertically by the cutting. Vertical apportioning is carried out by gathering traits into segments focused around the associations among the characteristics. Level apportioning is carried out by gathering tuple's into containers. Cutting jam

utility in light of the fact that it bunches profoundly corresponded traits together and jam the relationships between such qualities. At the point when the information set holds Qis and one SA, bucketization need to break their association. Cutting can assemble some QI traits with the SA for saving quality associations with the touchy characteristic. We display a novel system called cutting for security protecting information distribute.

II. RELATED WORK

Data Collection and Data Publishing: A typical scenario of data collection and publishing is described. In the data collection phase the data holder collects data from record owners. As shown in the fig.1 data-publishing phase the data holder releases the collected data to a data miner or the public who will then conduct data mining on the published data.

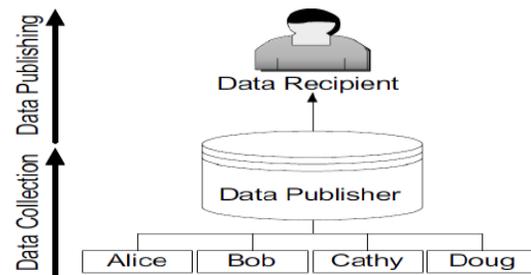


Figure 1: Data collection and Data Publishing

Security Preserving Data Publishing: The protection saving information distributed has the most fundamental structure that information holder has a table of the structure: D (Explicit Identifier, Quasi Identifier, Sensitive Attributes, non-Sensitive Attributes) holding data that expressly distinguishes record managers. Semi Identifier is a situated of characteristics that could conceivably distinguish

record managers. Touchy Attributes comprise of delicate individual particular data. Non-Sensitive Attributes holds all traits that don't fall into the past three classes.

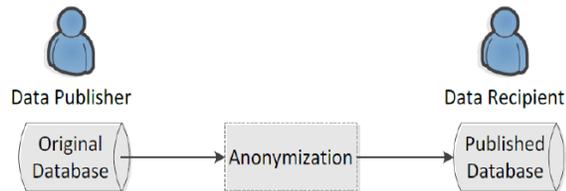


Figure 2: A Simple Model of PPDP

Information Anonymization: Data Anonymization is an innovation that changes over clear content into a non-intelligible structure. The procedure for protection saving information distributed has gotten a considerable measure of consideration lately. Most well known anonymization systems are Generalization and Bucketization. The fundamental distinction between the two-anonymization systems lies in that bucketization does not sum up the QI properties.

Generalization: Generalization is one of the ordinarily anonymized methodologies that supplant semi identifier values with values that are less particular yet semantically steady. All semi identifier values in a gathering would be summed up to the whole gathering degree in the QID space. In the event that no less than two transactions in a gathering have unique values in a certain segment then all data about that thing in the current gathering is lost. QID utilized as a part of this methodology incorporates all conceivable things in the log. With the goal generalization should be compelling, records in the same basin must be near one another so that summing up the records would not lose an excess of data. The information examiner need to make the uniform dispersion supposition that each quality in a summed

up interim/set is just as conceivable to perform information investigation or information mining assignments on the summed up table. This essentially decreases the information utility of the summed up information.

Bucketization: Bucketization is to segment the tuple's in T into cans and afterward to discrete the touchy property from the non-delicate ones by arbitrarily permuting the delicate trait values inside each one basin.

We utilize bucketization as the strategy for building the distributed information from the first table T. We apply an autonomous irregular stage to the segment holding S-values inside each one container. The ensuing set of basins is then distributed. While bucketization has preferred information utility over generalization it has a few limits. Bucketization does not anticipate enrollment revelation in light of the fact that bucketization distributes the QI values in their unique structures. Bucketization obliges an agreeable partition in the middle of Qis and Sas. In numerous information sets it is misty which properties are Qis and which are Sas. By differentiating the touchy characteristic from the QI properties. Bucketization breaks the quality relationships between the Qis and the Sas. The anonymized information comprise of a set of pails with permuted delicate quality qualities. Bucketization has been utilized for anonymizing high-dimensional information.

III. BASIC IDEA REGARDING SLICING

DATA SLICING method partitions the data both horizontally and vertically, which we discussed previously. The method partitions the data both

horizontally and vertically. This reduces the dimensionality of the data and preserves better data utility than bucketization and generalization.

Data slicing method consists of four stages:

o **Partitioning attributes and columns:** An attribute partition consists of several subsets of A that each attribute belongs to exactly one subset. Consider only one sensitive attribute S one can either consider them separately or consider their joint distribution.

o **Partitioning tuple's and buckets:** Each tuple belongs to exactly one subset and the subset of tuple's is called a bucket.

o **Generalization of buckets:** A column generalization maps each value to the region in which the value is contained.

o **Matching the buckets:** We have to check whether the buckets are matching.

Data Slicing: The original micro-data consist of quasi-identifying values and sensitive attributes. As shown in the fig.1 patient data in a hospital. Data consists of Age, Sex, Zip code, disease. A generalized table replaces values.

Age	Sex	Zip code	Disease
22	M	47906	Cancer
22	F	47906	Thyroid
33	F	47905	Thyroid
52	F	47905	Diabetes
54	M	47902	Thyroid
60	M	47902	Cancer
60	F	47904	Cancer

Table.1: Original microdata published.

The recoding that preserves the most information is “local recoding”. The first tuple are grouped into buckets and then for each bucket because same attribute value may be generalized differently when they appear in different buckets.

Age	Sex	Zip code	Disease
[20-52]	*	4790*	Cancer
[20-52]	*	4790*	Thyroid
[20-52]	*	4790*	Thyroid
[20-52]	*	4790*	Diabetes
[54-64]	*	4790*	Thyroid
[54-64]	*	4790*	Cancer
[54-64]	*	4790*	Cancer

Table.2: Generalized data

Table.2 shows the generalized data of the considered data in the above table. One column contains QI values and the other column contains SA values in bucketization also attributes are partitioned into columns. In the table.3 we describe the bucketization data. One separates the QI and SA values by randomly permuting the SA values in each bucket.

Age	Sex	Zip code	Disease
22	M	47906	Cancer
22	F	47906	Thyroid
33	F	47905	Thyroid
52	F	47905	Diabetes
54	M	47902	Thyroid
60	M	47902	Cancer
60	F	47904	Cancer

Table.3: Bucketized data

The basic idea of slicing is to break the association cross columns, to preserve the association within each column. It reduces the dimensionality of

data and preserves better utility. Data slicing can also handle high-dimensional data.

(Age, Sex)	(Zip code, Disease)
(22, M)	(47906, Cancer)
(22, F)	(47906, Thyroid)
(33, F)	(47905, Thyroid)
(52, F)	(47905, Diabetes)
(54, M)	(47902, Thyroid)
(60, M)	(47902, Cancer)
(60, F)	(47902, Cancer)

Table.4: Sliced data

IV. BACKGROUND APPROACH

Microdata publishing enable researchers and policy-makers to analyze the data and learn important information. Privacy is a key parameter in sensitive attribute disclosures. For privacy in Microdata publishing generalization and bucketization techniques based on k-anonymity, l-diversity approaches were used. Generalization fails to handle high dimensional data Bucketization fails to maintain clear separation between quasi-identifying attributes and sensitive attributes. K-anonymity protects against identity disclosures, but it does not provide sufficient protection against attribute disclosures. L-diversity protects against attribute disclosures but fails to prevent probabilistic attacks. So a better system is required that can with stand these failures and offers significant performance rise. For privacy in Microdata publishing a novel technique called slicing is used, which partitions the data both horizontally and vertically. Slicing preserves better data utility than generalization and can be used for membership

disclosure protection. Slicing can handle high-dimensional data. Attribute Partition and Columns

- a. Tuple Partition and Buckets
- b. Slicing
- c. Column Generalization

These methods compromise on overall data utility to maintain diversity requirement. A better system is required that can that can with stand high-dimensional data handling and sensitive attribute disclosure failures. Fig.3 describes the slicing architecture.

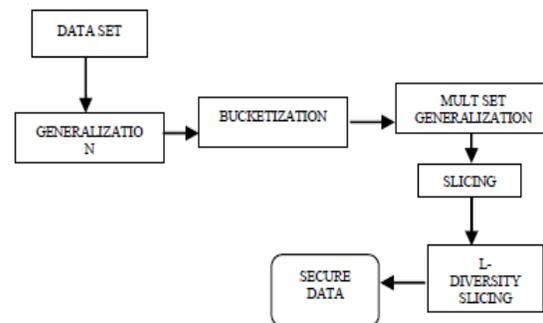


Figure 3: Slicing Architecture.

For Sliced data to obey the diversity requirement random grouping methods were used. Slicing algorithm consists of three phases: attribute partitioning, column generalization, and tuple partitioning. Involves the following procedures to attain data anonymity.

V. PROPOSED APPROACH

For privacy in Microdata publishing we still use slicing, which partitions the data both horizontally and vertically. Existing Slicing methods compromise on overall data utility to maintain diversity

requirement. Therefore, we propose to replace random grouping with more effective tuple grouping algorithms such as Tuple Space Search algorithm based on hashing techniques. A tuple is defined as a vector of k lengths, where k is the number of fields in a filter. For example, in a 5-field filter set, the tuple [7, 12, 8, 0, 16] means the length of the source IP address prefix is 7, the length of the destination IP address prefix is 12, the length of the protocol prefix is 8 (an exact protocol value), the length of the source port prefix is 0 (wildcard or "don't care"), and the length of the destination port prefix is 16 (an exact port value). We can partition the filters in a filter set to the different tuple groups. Since the filters in a same tuple group have the same tuple specification, they are mutual exclusive and none of them overlaps with others in this tuple group. Now we can perform the packet classification across all the tuple has to find the best-matched filter. If multiple tuple groups report matches, we resolve the best-matched filter by comparing their priorities. The filters in a tuple can be easily organized into a hash table, where we use the tuple specification to extract the proper number of bits from each field as the hash key. This key can be used for faster indexing, sorting and a primarily for accurate comparisons. The efficiency of tuple grouping algorithms enables its application to handle slicing problems that were previously prohibitive due to high-dimensional data handling and sensitive attribute disclosures.

Slicing With Tuple Grouping Algorithm:

Slicing with Tuple grouping algorithm provides efficient random tuple grouping for micro data publishing. In each column contains sliced bucket (SB) that permutated random values for each

partitioned data. The frequency of the value in each one of the scan's-diversity algorithm checks the diversity when the each sliced table

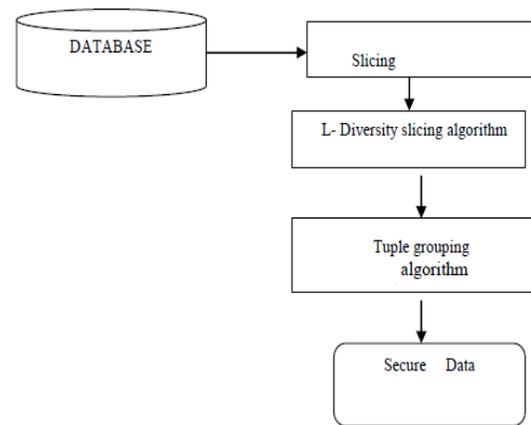


Figure 4: Architecture of slicing with tuple grouping

Step 1: Extract the data set from the database.

Step 2: Removes the queue of buckets and splits the Bucket into two

Step 3: computes the sliced table.

Step 4: Diversity maintains the multiple matching Buckets.

Step 5: Random tuple's are computed.

Figure 5: shows the algorithm that the tuple algorithm describes the functional procedure with respective to the architecture of the slicing with the tuple algorithm.

The main part of the tuple-partition algorithm is to check whether a sliced table satisfies „l-diversity gives a description of the diversity-check algorithm. The algorithm maintains a list of statistics $L(t)$ about t 's matching buckets. In each element in the list $L(t)$ contains statistics about one matching bucket b . The algorithm first takes one scan of each bucket b to record the frequency $f(v)$ of each column

value v in bucket b . The algorithm takes one scan of each tuple t in the table t to find out all tuple's that match b and record their matching probability $p(t, B)$ and the distribution of candidate sensitive values $d(t, B)$ which are added to the list $l(t)$. A final scan of the tuple's in t will compute the $p(t, b)$ values based on the law of total probability.

VI. EXPERIMENTAL EVALUATION

To allow direct comparison, we use the L -diversity for two anonymization techniques: slicing and optimized slicing for tuple grouping. We demonstrate experiment demonstrates that:

- Slicing preserves better data utility than generalization
- Slicing is more effective than bucketization in workloads involving the sensitive attribute
- The sliced table can be computed efficiently

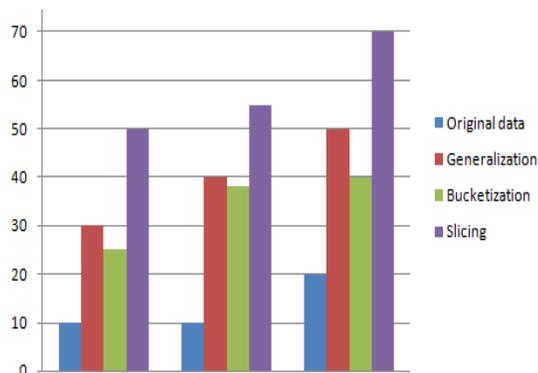


Figure 6: Computational efficiency

We compare slicing with optimized slicing in terms of computational efficiency. Fig.6 shows the computational efficiency.

VII. CONCLUSION

Slicing overcomes the limitations of generalization and bucketization and preserves better utility while protecting against privacy threats. That slicing preserves better data utility than generalization and is more effective than bucketization in workloads involving the sensitive attribute. Initially, we consider slicing where each attribute is in exactly one column. Our experiments show that random grouping is not very effective. Proposed grouping algorithm is optimized L -diversity slicing check algorithm obtains the more effective tuple grouping and Provides secure data. Data Slicing overcomes the limitations of generalization and bucketization and preserves better utility while protecting against privacy threats. Another important advantage of slicing is that it can handle high-dimensional data.

VIII. REFERENCES

- [1] Tiancheng Li, Ninghui Li, Jian Zhang, Ian Molloy, "Slicing: A New Approach for Privacy Preserving Data Publishing," IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 24, NO. 3, PP:561-574 ,MARCH 2012.
- [2] R.Maheswari, V.Gayathri, S.Jaya Prakash, " Tuple Grouping Strategy for Privacy Preservation of Microdata Disclosure," Proc. Int'l Conf. Very Large Data Bases (VLDB), pp. 901-909, 2005.
- [3] Amar Paul Singh, Ms. Dhanshri Parihar, " A Review of Privacy Preserving Data Publishing Technique," International Journal of Emerging

Research in Management &Technology, pp. 32-38, 2013.

[4] M.Alphonsa, V.Anandam, D.Baswaraj, "Methodology of Privacy Preserving Data Publishing by Data Slicing," INTERNATIONAL JOURNAL OF COMPUTER SCIENCE AND MOBILE APPLICATIONS, pp. 30-34, 2013.

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

[6] I. Dinur and K. Nissim, "Revealing Information while Preserving Privacy," Proc. ACM Symp. Principles of Database Systems (PODS), pp. 202-210, 2003.

[7] C. Dwork, "Differential Privacy," Proc. Int'l Colloquium Automata, Languages and Programming (ICALP), pp. 1-12, 2006.

[8] C. Dwork, "Differential Privacy: A Survey of Results," Proc. Fifth Int'l Conf. Theory and Applications of Models of Computation (TAMC), pp. 1-19, 2008.