

International Journal of Advance Research in Computer Science and Management Studies

Research Article / Survey Paper / Case Study

Available online at: www.ijarcsms.com

Anonymization Techniques for Preserving Highly Sensitive Information

Divya Naidu¹

Central College of Engineering and Management
Dept. of Computer Science and Engineering
Raipur, Chhattisgarh, India

Vaibhav Chandrakar²

Assistant Prof. Central College of Engineering and Management
Dept. of Computer Science and Engineering
Raipur, Chhattisgarh, India

Abstract: Securing data insurance is a basic issue in microdata circulation. Namelessness methodologies normally mean to guarantee solitary security, with immaterial impact on the idea of the released data. Starting late, a few models are familiar with ensure the security guaranteeing or conceivably to lessen the information setback to such a degree as could be permitted. That is, they moreover improve the flexibility of the mysterious framework to make it even more close to the real world, and after that to meet the different needs of everybody. Distinctive recommendation and counts have been planned for them meanwhile. Right now, presents two anonymization strategies for including vertices and edges for concealing valuable data. The calculation for including edges executes quicker when contrasted with calculation for including vertices however concealing data ability is more in vertices including calculation since it includes vertices inferred that it additionally include edges so it turns out to be increasingly strong.

Keywords: Anonymity techniques, anonymity models, privacy preserving algorithm.

I. INTRODUCTION

The present databases contain an extensive proportion of fragile individual data. So it's huge to layout information systems which may restrict the important of individual data. For example, consider a recuperating focus that keeps up tolerant records [1] [2].

The recuperating office longings to reveal data to an association in such some manner that the association can't find that patients have that ailments. One framework to officially show insurance courses of action is to specific touchy data as requests and actualizes phenomenal security, an appallingly solid idea of security that guarantees that the other inquiry answered by the information won't reveal any data regarding the sensitive data [3] [4].

A. Security Preserving Data Publishing

Individual records of people are sensibly being assembled by different government and association establishments for the necessities of data assessment [7] [8]. The data assessment is urged by these relationships to circulate "sufficiently private" contemplations over this information that are gathered. Privacy could be a twofold edged brand - there ought to be adequate insurance to guarantee that tricky data concerning everybody isn't uncovered by the viewpoints and at a similar time there ought to be adequate information to play out the examination. Additionally, a foe who needs to assemble fragile data from the revealed points of view now and again has a few information concerning everyone inside the information. The rule objective is to change over the primary information into some puzzling sort to keep from prompting its record proprietor's delicate data as inspected in [9].

B. Information Anonymization

Data anonymization is the path toward removing before long recognizable information from educational files, to make everybody obscure about whom the data portray. It permits the trading of data over a breaking point, as between two workplaces inside concentration or between two workplaces, however diminishing the danger of unplanned revealing, and in bound conditions in an exceedingly way that awards examination and assessment post-anonymization [10] [11]. This framework is used as a piece of endeavors to extend the security of the data while empowering the data to be destitute down and used. It changes the data that will be used or disseminated to keep the distinctive evidence of key information. Data anonymization strategies, for instance, k-secrecy, l-decent variety characteristics what's more, t-closeness are expansive.

k-Anonymity: The basic course of action of k-obscurity is to shield a dataset against re-recognized by summarizing the attributes that might be utilized in a linkage assaults (semi identifiers). An informational collection is viewed as k-mysterious if every datum thing can't be perceived from at any rate k-1 elective data things [12].

l-Diversity: l-decent variety characteristics could be a combination of gathering based generally anonymization that is wont to defend security in learning sets by diminishing the coarseness of a learning depiction. This decreasing may be a trade off that breezes up in some loss of amplexness of data organization or mining calculation so as to achieve some security. The l-varying characteristics model is connected degree development of the k-mystery show that decreases the cruelty of information portrayal misuse systems and what's more hypothesis and disguise demonstrated any given record maps onto in any event k elective records inside the data [13].

t-Closeness t-closeness could be an additional refinement of l-grouped characteristics pack based for the most part anonymization that is accustomed shield security in learning sets by reducing the coarseness of an information depiction. t-closeness could be an extra refinement of l-arranged characteristics bunch basically based anonymization that is wont to spare insurance in learning sets by diminishing the coarseness of an information outline. This diminishing could be a trade off that breezes up in some loss of practicality of data organization or mining calculation so as to comprehend a couple of security [14].

II. K-ANONYMITY

k-Anonymity could be a conventional model of security [16]. The goal is to outline each record hazy from a delineated assortment (k) records if attempts district unit made to distinguish the information. A course of action of information is k-anonymized if, for any record with a given plan of qualities, there square measure in any occasion k-1 elective records that coordinate these characteristics. The properties can be any of the going with sorts.

The utilization of k-obscurity needs the preliminary ID of the semi identifier. The semi identifier relies upon the external data available to the recipient, since it chooses the associating limit (not all possible outside tables extend unit open to each potential learning recipient); and assorted semi identifiers will point of fact exist for a given table [15].

Example:

If the recently referenced table is to be anonymized with Anonymization Level (AL) set to 2 and the course of action of Quasi identifiers as $QI = \{AGE, SEX, ZIP, PHONE\}$. Sensitive quality = $\{SALARY\}$. The semi identifiers and sensitive characteristics are recognized by the relationship as showed by their principles and guideline.

TABLE I: Table to be Anonymized

ID	Age	Sex	Zip	Phone	Salary (in Rs.)
1	24	M	641015	9994258665	78000
2	23	F	641254	9994158624	45000
3	45	M	610002	8975864121	85000
4	34	M	623410	7456812312	20000

TABLE II: Anonymized Table

ID	Age	Sex	Zip	Phone	Salary (in Rs.)
*	20-50	ANY	641***	999*****	78000
*	20-50	ANY	641***	999*****	45000
*	20-50	ANY	612***	897*****	85000
*	20-50	ANY	623***	745*****	20000

A. Generalization

Speculation is the path toward changing over a motivating force into a less specific general term. For ex, "Male" and "Female" can be summed up to "Any". At the going with levels speculation techniques can be associated.

1. Attribute (AG): Generalization is performed at the fragment level; all the characteristics in the segment are summed up at a hypothesis step.
2. Cell (CG): Generalization can similarly be performed on a singular cell; finally a summarized table may contain, for a specific area and qualities at different degrees of speculation.

B. Suppression

Concealment involves in turning away sensitive data by emptying it. Concealment can be associated at the degree of single cell, entire tuple, or entire fragment, licenses lessening the proportion of theory to be compelled to achieve k-obscurity.

1. Tuple (TS): Suppression is performed at section level; concealment activity clears whole tuple.
2. Attribute (AS): Suppression is performed at portion level; concealment activity covers all of the estimations of a fragment.

III. LITERATURE SURVEY

Xuyun Zhang et al. [16] proposes giving security and insurance over the middle of the road informational collections become contest issue since enemies may hold small scale information by recognizing various information records. Encryption of all datasets when all is said in done society arrange called cloud take in past frameworks may very dull and extravagant.

Mohammad Reza Zare Mirakabad et al. [17] focuses giving insurance over the data creation. Under security data use and abhorrence of disclosure of individual character is increasingly basic. One of the data anonymization strategies called K-mystery keeps the exposure of individual character anyway it is generally fail to achieve.

Min Wu et al. [18] proposes sparing security is most fundamental anyway a comparable time it is bother in appearance of little scope data release. In the point of view of characteristic revelation K-anonymity isn't well. So, we propose new framework called an ordinal division-based affectability mind full contrasting characteristics metric model.

Yunli Wang et al. [19] proposes k-secrecy fails to achieve characteristics disclosure anyway in l-grouped characteristics intends to achieve trademark introduction. Second data anonymization technique centre around cutting the illation from freed scaled down scale attributes.

Jordi Soria Comas et al. [20] focuses data anonymization methodologies spare insurance, k-obscurity and €-differential security are two standard assurance show. The t-closeness is the enlargement of k-indefinite quality, the improvement of private touchy information relies upon Bucketization calculation.

IV. METHODOLOGY

In this section we present the two Anonymization techniques

- a. With Vertices Anonymization

b. With Edges Anonymization

A. With Vertices Anonymization

In this the vertices are anonymized. The additional vertices are deliberately added to the network or graph to hide the degree of information present in it.

B. With Edges Anonymization

In this the edges are anonymized. The additional edges are deliberately added to the network or graph to hide the degree of information present in it.

Adding vertices and edges are depended upon the k-anonymization algorithm.

Algorithm: Edges Anonymization
Input: An initial multi sensitive graph $G(V, E)$
Output: Graph $G'(V', E')$ – with added edges
<ol style="list-style-type: none"> 1. Get degrees in descending order 2. Anonymize the degree 3. Using anonymize vector, add additional degree 4. Create subgraph for the degree vector 5. Return Anonymized graph

Fig. 2. Shows the algorithm of Anonymization using Edges

For Vertex Anonymization – Vertex Anonymization algorithm is used. The algorithm effectively anonymized the data. The algorithm is presented in fig. 2.

Algorithm: Vertices Anonymization
Input: An initial multi sensitive graph $G(V, E)$
Output: Graph $G'(V', E')$ – with added vertices
<ol style="list-style-type: none"> 1. fetch orbits from the graph using stab graph Algorithm 2. Iterate through the orbits of the graph <ol style="list-style-type: none"> a. Introduce new vertex and add to the graph and include it into orbit b. Get ID of the vertex c. Connect new edges according to the orbit d. In same orbit connect them by tag and in different connect them by the regular graph connection 3. Return Anonymized graph

For adding least number of edges in the graph the following equation is used:

If

$$L_1(\hat{\mathbf{d}} - \mathbf{d}) = \sum_i |\hat{\mathbf{d}}(i) - \mathbf{d}(i)|,$$

$$GA(\hat{G}, G) = |\hat{E}| - |E| = \frac{1}{2} L_1(\hat{\mathbf{d}} - \mathbf{d}).$$

Where G is the Graph with E edges and V vertices.

V. RESULTS

The examination are led utilizing Eclipse structure on java language. We have shown two calculation for anonymization. Right off the bat by including edges and besides by including vertices. As including vertices is unpredictable procedure since we can't just include vertices into the diagram, we need likewise to include edges naturally. Consequently, time taken by the including vertices is long when contrasted with including edges. Fig. 4. Shows the subtleties of Facebook organize dataset.

Degree ^	Vertices		
1	75	Total Vertices	4039
2	98	Vertices added (%)	0 (0.00%)
3	93	Total Edges	88234
4	99	Edges added (%)	0 (0.00%)
5	93	Duration	0.00sec
6	98		
7	98		
8	111		

Fig. 4. Snapshot of Facebook circle dataset

Table: IV. 5-Anonymized – Adding Edges Output

Total Vertices	4039
Vertices Added	0
Total Edges	95420
Edges Added	7186 (8.14%)
Time Taken	3.59 sec

Table: IV. 5-Anonymized – Adding Vertices Output

Total Vertices	4386
Vertices Added	347 (8.59%)
Total Edges	181487
Edges Added	93253 (105.69%)
Time Taken	16.65 sec

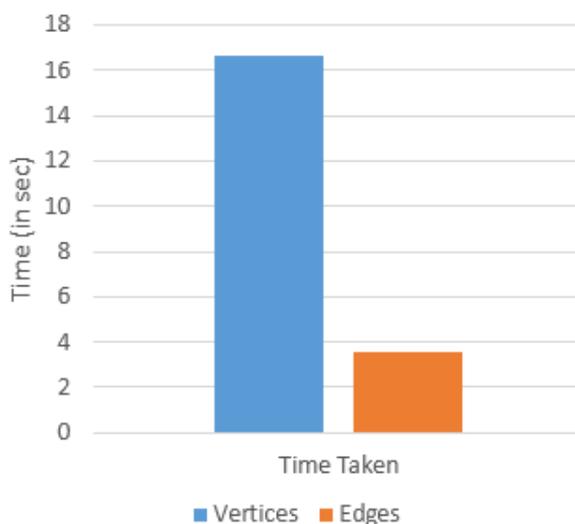


Fig. 5. Shows the time taken for execution of two algorithms

VI. CONCLUSION

The level of a center point in a chart arrange, among other helper characteristics, can to an immense degree perceive the center point from various centers. Right now, focused a particular graph anonymity believed that keeps the re-distinguishing proof of individuals by an attacker with sure prior data of the degrees. We officially characterized the Graph Anonymization issue that, given a data chart demands the base number of edge augmentations (or deletions) that empower the difference in the commitment to a degree-mysterious chart i.e., an outline wherein every center point has a comparable degree with k-1 distinct centers.

On the off chance that by change programmer knows the specific data, it can't cross the entire arrangement of data. This is the idea of anonymization. Consequently, it is utilized by numerous associations, for example, Hospitals, Collages, Universities, Secret Data holder to conceal basic data.

References

1. M. E. Kabir, H. Wang and E. Bertino, "Efficient systematic clustering method for k-anonymization," *Acta Informatica*, Springer, Vol. 48, 2011, pp. 51-66.
2. J. W. Byun, A. Kamra, E. Bertino, and N. Li, "Efficient k-anonymization using clustering techniques," in *Proceedings of International Conference on Database Systems for Advanced Applications*, 2007, pp. 188-200.
3. X. Xiao and Y. Tao, "Anatomy: simple and effective privacy preservation," in *Proceedings of the 32nd International Conference on Very Large Data Bases*, 2006, pp.139-150.
4. Xuyun Zhang, Chang Liu, Surya Nepal, Chi Yang, Wanchun Dou, Jinjun Chen" Combining Top-Down and Bottom-Up: Scalable Sub-Tree anonymization over Big data using MapReduce on Cloud".
5. J. Goldberger and T. Tassa, "Efficient anonymization with enhanced utility," *Transactions on Data Privacy*, Vol. 3, 2010, pp. 149-175.
6. M. Terrovitis, N. Mamoulis, and P. Kalnis. "Privacy-preserving anonymization of set-valued data." *PVLDB*, 1(1):115–125, 2008.
7. Md Nurul Huda, Shigeki Yamada, and Noboru Sonehara, "On Enhancing Utility in k-Anonymization", *International Journal of Computer Theory and Engineering*, Vol. 4, No. 4, August 2012.
8. Pawan R. Bhaladhare and Devesh C. Jinwala, "Novel Approaches for Privacy Preserving Data Mining in k-Anonymity Model" , *JOURNAL OF INFORMATION SCIENCE AND ENGINEERING* 32, 63-78 (2016).
9. Mohammed, N. and Fung, B. C. M, "Centralized and distributed anonymization for high-dimensional healthcare data", *ACM Trans. Knowl. Discov. Data.* 4, 4, Article 18 (October 2010), 33 pages.
10. S. E. Fienberg, A. Slavkovic and C. Uhler," Privacy Preserving GWAS Data Sharing", 2011 11th IEEE International Conference on Data Mining Workshops.
11. A. G. Divanis and G. Loukides," PCTA: Privacy-constrained Clustering-based Transaction
12. Data Anonymization", *ACM* 2011.
13. S. Kisilevich, L. Rokach, Y. Elovici, and B. Shapira. "Efficient multidimensional suppression for k-anonymity." *TKDE*, 22:334–347, 2010.
14. G. Loukides, A. Gkoulalas-Divanis, and B. Malin. Anonymization of electronic medical records for validating genome-wide association studies. *PNAS*, 17:7898–7903, 2010.
15. J. Cao, P. Karras, C. Ra'issi, and K. Tan. rho-uncertainty: Inference-proof transaction anonymization. *PVLDB*, 3(1):1033–1044, 2010.
16. Loukides, G., Shao, J.: Capturing data usefulness and privacy protection in k-anonymisation. In: *Proceedings of the 2007 ACM Symposium on Applied Computing* (2007)
17. X. Zhang, C. Liu, S. Nepal, S. Pandey and J. Chen, "A Privacy Leakage Upper Bound Constraint-Based Approach for Cost-Effective Privacy Preserving of Intermediate Data Sets in Cloud," *IEEE Transactions on Parallel and Distributed Systems*, vol. 24, no. 6, pp. 1192-1202, 2013.
18. Mohammad Reza Zare Mirakabad School of Computer Sciences, USM, Malaysia Intern at School of Computing, NUS, Singapore reza@cs.usm.my, reza.z@comp.nus.edu.sg "Diversity versus Anonymity for Privacy Preservation".
19. Min Wu, Xiaojun Ye Institute of Information System and Engineering School of Software, Tsinghua University, Beijing, 100084, China "Towards the Diversity of Sensitive Attributes in k-Anonymity".
20. Yunli Wang, Yan Cui, Liqiang Geng and Hongyu Liu, "A new perspective of privacy protection: Unique distinct l-SR diversity," 2010 Eighth International Conference on Privacy, Security and Trust, Ottawa, ON, 2010.