

A Survey on Anonymization Based Privacy Preserving Models

Teshu Kant Parkar

M.Tech Scholar, Computer Science and Engineering, Shri Shankaracharya Group of Institution, Bhilai (C.G.), India

Yamini Chouhan

Assistant Professor, Department of Computer Science and Engineering, Shri Shankaracharya Group of Institution, Bhilai (C.G.), India

Samta Gajbhiye

Professor, Head of Department, Computer Science and Engineering, Shri Shankaracharya Group of Institution, Bhilai (C.G.), India

Abstract – An expanding number of associations keep up collection of information about people. Hospitals keep medical records of their patients, business organizations gather data of their customers, and web benefit organizations monitor the inclinations of their clients. Publication of this information can be valuable for investigating, epidemic examinations, business advancement, measurable investigation, and so forth. Information appropriation could be either open, as a publication on the web, or confined, for instance to specific research groups. In either case, the information holder must guarantee the privacy of people whose individual data is incorporated into the released dataset. In this paper, a survey for existing strategies and to analyze the strength and limitations of these methodologies is talked about.

Index Terms – Anonymity techniques, anonymity models, privacy preserving, algorithm.

1. INTRODUCTION

Privacy is a standout amongst the most concerned issues in cloud computing. Personal information like monetary exchange records and electronic health records are to a great degree touchy in spite of the fact that that can be dissected and mined by inquire about association. Information privacy issues should be tended to before informational collections are shared on cloud for examination reason. Information anonymization alludes to as concealing delicate information for proprietors of information records [1].

Substantial scale informational indexes are generalized utilizing two stage top-down specialization for information anonymization. This procedure split into two stages. In the principal phase, original informational collections are parceled into a gathering of littler informational collections, and these informational indexes are anonymized in parallel, producing middle of the road comes about. In the second stage, the middle of the road comes about are coordinated into one, and further anonymized to accomplish predictable k-mysterious

[6] informational collections. In such cases, informational collections are anonymized as opposed to encrypted to guarantee the two information utility and privacy preserving. Instead of encryption, anonymization of information should be possible as a methods for preserving privacy [2].

Information anonymization is generally utilized strategy for Privacy Preserving of information in non-intelligent information distributing situation Data anonymization alludes to the concealing the personality or delicate information for proprietors information record. The privacy of individual can be adequately safeguarded while some total data is shared for information investigation and mining. Several models of security can enhance Data Anonymization incorporate k-anonymity and l-diversity.

2. PRIVACY PRESERVING APPROACHES

A. K-Anonymity

K-anonymity is a property controlled by certain anonymized data. Given individual particular field structured information; create an arrival of the information with logical ensures that the people who are the subjects of the information can't be re-recognized while the information remain basically valuable. An arrival of information is said to have the k-anonymity property if the data for every individual contained in the discharge can't be recognized from at any rate k-1 people.

B. K-Anonymization Methods

Suppression:

In this technique, certain estimations of the attributes are supplanted by a mark '*'. All or a few estimations of a segment might be supplanted by '*'. Suppression comprises in averting delicate data by evaluating it. Suppression can be connected at the level of single cell, whole tuple, or whole segment, permits

diminishing the measure of speculation to be forced to accomplish k-anonymity.

- 1) Tuple (TS): Suppression is performed at column level; suppression operation evacuates entire tuple
- 2) Attribute (AS): Suppression is performed at segment level; suppression operation shrouds every one of the estimations of a segment.
- 3) Cell (CS): Suppression is performed at single cell level; at long last k-anonymized table may wipe out just certain cells of a given tuple/quality.

Generalization:

In this technique, singular estimations of attributes are supplanted by with a more extensive class. For instance, the value '19' of the attribute 'Age' might be replaced by ' ≤ 20 ', the value '23' by ' $20 < \text{Age} \leq 30$ '. Generalization is the way toward changing over an incentive into a less particular general term. For ex, "Male" and "Female" can be generalized to "Any". At the distinguish levels generalization procedures can be connected.

- 1) Attribute (AG): Generalization is performed at the segment level; all the qualities in the section are generalized at a speculation step.
- 2) Cell (CG): Generalization can likewise be performed on a single cell; at long last a summed up table may carry, for a particular section and values at different levels of generalization.

L-Diversity:

L-diversity is a type of gathering based anonymization that is utilized to save .The l-diversity display is an expansion of the k-anonymity show which decreases the granularity of information portrayal utilizing methods including speculation and suppression to such an extent that any given record maps onto at any rate k different records in the information.

The l-diversity display handles a portion of the shortcomings in the k-anonymity show where ensured characters to the level of k-individual isn't identical to securing the comparing touchy esteems that were generalized or stifled, particularly when the delicate esteems inside a gathering display homogeneity.

T-Closeness

Given the presence of assaults where touchy attributes might be construed based upon the circulation of qualities for l-diversity information, the t-closeness technique was made to advance l-diversity by also keeping up the dispersion of delicate fields.

An identicalness class is said to have t-closeness if the separation between the appropriation of a delicate attribute in

this class and the dissemination of the attribute in the entire table is close to an edge t.

Bottom-Up Generalization

By imaginatively applying Map Reduce on cloud to Bottom Up Generalization (BUG) for information anonymization and purposely outline a gathering of inventive Map Reduce employments to solidly achieve the generalizations in an exceedingly adaptable manner. Besides, present a versatile Advanced BUG approach, which performs speculation on various parceled informational index and the subsequent transitional anonymizations are converged to discover last anonymization which is utilized to anonymize the first informational collection. Results demonstrate that our approach can fundamentally enhance the versatility and effectiveness of BUG for information anonymization over existing methodologies.

Top-Down Specialization

By and large, TDS is an iterative procedure beginning from the topmost area esteems in the scientific classification trees of attributes. Each round of emphasis comprises of three stages [4]:

- Searching the best specialization
- Performing specialization
- Refreshing values of the search metric for the next round

Such a process is repeated until the point when k-anonymity is violated, to uncover the most extreme information utility. The decency of a specialization is estimated by an inquiry metric

3. RELATED WORK

M. E. Kabir et al. [1], This paper introduces a clustering (Clustering partition record into clusters to such an extent that records inside a cluster are like each other, while records in various clusters are most unmistakable from each other.) based k-anonymization method to limit the data misfortune while in the meantime guaranteeing information quality. Privacy preservation of people has attracted significant interest's information mining research. Anonymization strategies by means of generalization or suppression can ensure private data, however lose esteemed data. The test is the way to limit the data misfortune amid the anonymization procedure. Author allude to the test as a methodical clustering issue for k-anonymization which is analyzed in this paper. The proposed method embraces cluster comparable information together and after that anonymizes each gathering exclusively. The structure of systematic clustering issue is characterized and researched through worldview and properties.

J. W. Byun et al. [2], K-anonymization strategies have been the concentration of serious research over the most recent

couple of years. An imperative necessity for such procedures is to guarantee anonymization of information while in the meantime limiting the data misfortune coming about because of information changes. In this paper, author propose an approach that uses clustering to limit data misfortune and in this way guarantee great information quality. The key perception here is that information records that are normally like each other ought to be a piece of a similar identicalness class. Author hence detail a particular clustering issue, alluded to as k-part clustering issue. Author demonstrate that this issue is NP-hard and exhibit an insatiable heuristic, the complexity of which is in $O(n^2)$. As a major aspect of our approach Author build up an appropriate metric to appraise the data misfortune presented by generalizations, which works for both numeric and absolute information.

X. Xiao et al. [3], this paper exhibits a novel procedure, life systems, for distributing delicate information. Life structures discharges all the semi identifier and delicate esteems specifically in two separate tables. Joined with a gathering component, this approach ensures protection, and catches a lot of connection in the microdata. We build up a direct time algorithm for processing dissected tables that comply with the l-decent variety protection necessity, and limit the blunder of reproducing the microdata. Broad investigations affirm that our strategy permits essentially more successful information examination than the regular publication technique based on speculation. In particular, life systems grants total prevailing upon normal blunder underneath 10%, which is lower than the mistake acquired from a generalized table by requests of size.

Xuyun Zhang et al. [4], In big information applications, information security is a standout amongst the most concerned issues since preparing expansive scale protection delicate informational collections regularly requires calculation control gave by public cloud administrations. Sub-tree information anonymization, accomplishing a decent exchange off between information utility and mutilation, is a broadly received plan to anonymize informational indexes for security protection. Top-Down Specialization (TDS) and Bottom-Up Generalization (BUG) are two approaches to satisfy sub-tree anonymization. Nonetheless, existing methodologies for sub-tree anonymization miss the mark regarding parallelization ability, in this way inadequate with regards to adaptability in dealing with enormous information on cloud.

J. Goldberger et al. [5], the k-anonymization strategy is an ordinarily utilized protection safeguarding procedure. Past investigations utilized different measures of utility that go for upgrading the connection between the first public information and the generalized public information. We, remembering that an essential objective in discharging the anonymized database for information mining is to reason strategies for foreseeing the private information from the public information, propose another data theoretic measure that goes for upgrading the

relationship between the generalized public information and the private information. Such a measure essentially upgrades the utility of the discharged anonymized database for information mining. Author at that point continue to depict another and exceedingly productive algorithm that is intended to accomplish k-obscurety with high utility. That algorithm is based on an adjusted rendition of successive clustering which is the strategy for decision in clustering, and it is autonomous of the fundamental measure of utility.

M. Terrovitis et al. [6], in this paper, we contemplate the issue of securing protection in the publication of set-esteemed information. Consider an accumulation of value-based information that contains definite data about things purchased together by people. Indeed, even subsequent to expelling every single individual normal for the purchaser, which can fill in as connections to his personality, the publication of such information is as yet subject to security assaults from foes who have halfway learning about the set. Dissimilar to most past works, Author don't recognize information as touchy and non-delicate, however we think of them as both as potential semi identifiers and potential touchy information, contingent upon the perspective of the foe. Author characterize another form of the k-anonymity ensure, the k m-anonymity, to restrain the impacts of the information dimensionality and we propose proficient algorithms to change the database. Our anonymization show depends on speculation rather than concealment, which is the most widely recognized practice in related deals with such information. We build up an algorithm which finds the ideal arrangement, notwithstanding, at a high cost which makes it inapplicable for extensive, reasonable issues.

Md Nurul Huda et al. [7] k-anonymity is a standout amongst the most contemplated models of security safeguarding innovation. It restricts the connecting certainty between particular delicate data and a particular individual by concealing the distinguishing pieces of proof of every person into in any event k-1 others in the database. A k-anonymization algorithm is normally assessed utilizing data misfortune or information utility measurements. In this paper, we initially propose another quality metric, called the Efficiency metric. This metric beats the constraints of existing one dimensional measurements, speaking to either protection measure or information utility measure, utilized as a part of security safeguarding information sharing. We at that point introduce another heuristic algorithm for k-anonymization that offers high information utility and in addition an abnormal state of security.

Pawan R et al. [8], in security protecting information mining, anonymization based methodologies have been utilized to save the security of a person. Existing writing tends to different anonymization based methodologies for safeguarding the touchy private data of a person. The k-anonymity display is one

of the broadly utilized anonymization based approach. Be that as it may, the anonymization based methodologies experience the ill effects of the issue of data misfortune. To limit the data misfortune different cutting edge anonymization based clustering approaches viz. Eager k-part algorithm and Systematic clustering algorithm have been proposed. Among them, the Systematic clustering algorithm gives lesser data misfortune. What's more, these methodologies make utilization of all properties amid the making of an anonymized database. Thusly, the danger of exposure of delicate private information is higher by means of publication of the considerable number of properties. In this paper, Author propose two methodologies for limiting the exposure hazard and saving the security by utilizing efficient clustering algorithm. In the first place approach makes an unequal blend of semi identifier and touchy characteristic. Second approach makes an equivalent mix of semi identifier and touchy property. Author likewise assess our approach observationally concentrating on the data misfortune and execution time as imperative measurements. We outline the viability of the proposed approaches by contrasting them and the current clustering algorithms.

Mohammed N. et al. [9], sharing healthcare information has turned into a key prerequisite in healthcare framework administration; be that as it may, improper sharing and utilization of healthcare information could debilitate patients'

protection. In this article, we think about the protection worries of sharing patient data between the Hong Kong Red Cross Blood Transfusion Service (BTS) and the public hospitals. Author sum up their data and security necessities to the issues of brought together anonymization and conveyed anonymization, and recognize the significant difficulties that make customary information anonymization strategies not pertinent. Besides, Author propose another protection display called LKC-security to conquer the difficulties and present two anonymization algorithms to accomplish LKC-protection in both the concentrated and the circulated situations.

S. E. Fienberg et al. [10], customary factual techniques for the classification assurance for measurable databases don't scale well to manage GWAS (far reaching affiliation contemplates) databases and outside data on them. The later idea of differential protection, presented by the cryptographic group, is an approach which furnishes a thorough meaning of security with important security ensures within the sight of discretionary outside data. Expanding on such ideas, Author propose new strategies to discharge total GWAS information without trading off a person's security. Author show strategies for discharging differentially private minor allele frequencies, chisquare insights and p-values. Author think about these methodologies on mimicked information and on a GWAS investigation of canine hair length including 685 puppies.

TABLE I. Comparisons of various techniques and method used in existing system

| Ref. No. | Method Used | Data Source | Approach | Strength | Limitation |
|----------|-----------------------------------|------------------------------|--|---|---|
| [1] | Systematic clustering method | Medical data | Author presents a clustering based <i>k</i> -anonymization technique to minimize the information loss while at the same time assuring data quality | Results show that our method attains a reasonable dominance with respect to both information loss and execution time. | Need to extend the systematic clustering algorithm to <i>k</i> -anonymity model |
| [2] | <i>k</i> -anonymization algorithm | <i>k</i> -anonymized dataset | Author proposed an efficient <i>k</i> -anonymization algorithm by transforming the <i>k</i> -anonymity problem to the <i>k</i> -member clustering problem. | Author develop a suitable metric to estimate the information loss introduced by generalizations, which works for both numeric and categorical data. | Additional performance measures are there. |

| | | | | | |
|-----|---------------------------------------|-----------------------|---|---|---|
| [3] | Nearly-Optimal Anatomizing Algorithm | CENSUS dataset | Author develop a linear-time algorithm for computing anatomized tables that obey the l-diversity privacy requirement, and minimize the error of reconstructing the microdata. | experiments confirm that anatomy permits highly accurate aggregate information about the unknown microdata, with an average error below 10% | Need extend the technique to multiple sensitive attributes is an interesting topic |
| [4] | heuristic algorithm | Real world dataset | Proposed an approach that identifies which part of intermediate data sets needs to be encrypted while the rest does not, in order to save the privacy preserving cost. | results demonstrate that the privacy-preserving cost of intermediate data sets can be significantly reduced | Need to investigate privacy aware efficient scheduling of intermediate data sets in cloud |
| [5] | Sequential clustering algorithm. | US Census Bureau | Proposed the private mutual information (PMI) utility measure that aims at maximizing the correlation between the generalized public data and the private data. | PMI measure is much more suitable when the goal is to achieve anonymizations | Additional performance measures are there. |
| [6] | apriori-based anonymization algorithm | Real world dataset | Author develop an algorithm which finds the optimal solution, however, at a high cost which makes it inapplicable for large, realistic problems. | proposed algorithms are experimentally evaluated using real datasets | Need to consider sensitive values associated to set-valued quasi-identifiers. |
| [7] | LowCost algorithm | anonymous data record | Propose the Efficiency metric \square that represents both the utility and privacy of the anonymous data and can assess | Performs the best both in terms of utility measure and privacy measure. | Additional performance measures are there. |

| | | | | | |
|------|--|--|---|---|---|
| | | | anonymization algorithms fairly. | | |
| [8] | Greedy k -member algorithm and Systematic clustering algorithm | UCI machine learning repository database | Author propose two approaches for minimizing the disclosure risk and preserving the privacy by using systematic clustering algorithm. | Author illustrate the effectiveness of the proposed approaches by comparing them with the existing clustering algorithms. | Need to evaluate the performance on a combination of multiple SA and QI attributes. |
| [9] | centralized and distributed anonymization algorithm | Health care data | PPDP Has Received A Great Deal Of Attention In The Database And Data Mining Research Communities | Degradation Of Data / Service Quality Loss Of Valuable Information Increased Costs | Additional performance measures are there. |
| [10] | C-Differentially Private Algorithm | GWAS data | Author present methods for releasing differentially private minor allele frequencies, chisquare statistics and p -values. | Author provided a differentially private algorithm for releasing these statistics for the most relevant SNPs | Improvement required to get more accurate result |

4. CONCLUSION

In this paper, we talked about the Privacy preserving approaches in data publishing. We likewise talked about various anonymization method and for the most part focused on k -anonymity which include both generalization and suppression. The last part is about the generalization algorithm and its execution for securing the protection of information utilized for the most part for data analysis. In this paper, the discussion about the anonymity models, the significant execution ways and the techniques of anonymity algorithm, and also analyzed their strength and limitations.

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.