

One Size Does Not Fit All: Trade-offs between Misuse Probability and Level of Sanitization for Big Data

D. Radhika, D. Aruna Kumari

Abstract - Big data privacy has assumed importance as the cloud computing became a phenomenal success in providing a remote platform for sharing computing resources without geographical and time restrictions. However, the privacy concerns on the big data being outsourced to public cloud storage are still exist. Different anonymity or sanitization techniques came into existence for protecting big data from privacy attacks. In our prior works, we have proposed a misusability probability based metric to know the probable percentage of misusability. We additionally planned a system that suggests level of sanitization before actually applying privacy protection to big data. It was based on misusability probability. In this paper, our focus is on further evaluation of our misuse probability based sanitization of big data approach by defining an algorithm which will analyse the trade-offs between misuse probability and level of sanitization. It throws light into the proposed framework and misusability measure besides evaluation of the framework with an empirical study. Empirical study is made in public cloud environment with Amazon EC2 (compute engine), S3 (storage service) and EMR (MapReduce framework). The experimental results revealed the dynamics of the trade-offs between them. The insights help in making well informed decisions while sanitizing big data to ensure that it is protected without losing utility required.

Index Terms: Big data, privacy of big data, sanitization, misuse probability, utility of big data

I. INTRODUCTION

With Internet based computing paradigm, computing resources are provided with scalability and availability. This has led to outsourcing data to public cloud. Moreover, the data has grown to assume big data characteristics. Thus the outsourced data is to be subjected to data analytics for business intelligence. Such data is prone to various kinds of inference attacks [2]. Therefore, there is need for protecting privacy of big data. Big data and its processing frameworks like Hadoop [1] are widely used for business solutions. In this context different algorithms in data mining came with MapReduce flavour of programming [1], [23]. As big data is widely used in the public domain which is untrusted, there is emphasis on the sanitization of big data for protecting its privacy [7], [12], [13]. There are many solutions found in the literature with respect to privacy of big data. However, there is little research on the controlled way of anonymizing data. In other words, anonymization without precautions may prompt loss of utility.

In general, the increase in anonymity leads to decrease in utility of big data when it is subjected to analytics. This is very important hypothesis that is not only proved in this paper but also the trade-off between misuse probability of big data and sanitization are analysed. To be more precise, in this paper we propose a misuse probability score measure (also found in our prior work [2]) that is the basis for determining appropriate sanitization which strikes balance between utility and privacy of big data.

Big data might have multiple dimensions. Each attribute may belong to different kind of identifiers like general, sensitive and quasi identifiers. General identifiers need no protection. Sensitive identifiers need to be sanitized while quasi identifiers need different level of sanitization based on their level of misuse. Therefore, a misuse probability based metric is proposed (inspired by the work of Harel *et al.* [24]). This measure plays vital role in the appropriate sanitization of big data. In fact, it takes dataset as input and finds its misuse probability score. This score is then used to determine the level of sanitization needed. This phenomenon leads to achieve privacy of big data and ensure that its utility is not lost unnecessarily due to sanitization. Our commitments in this paper are as per the following:

- 1 We Projected a System for misuse probability score based sanitization of big data for protecting privacy.
- 2 We proposed an algorithm named Misuse Probability Based Optimal Big Data Sanitization (MP-OBS) to achieve the objective of finding trade-off between misuse probability and level of sanitization which is linked to the trade-off between sanitization and utility of big data. It helps in optimal privacy protection to big data.
- 3 We implemented the framework with cloud environment that includes Amazon EC2, Amazon EMR and Amazon S3. The framework is evaluated with standard measures like precision, recall, RMSE and classification

Accuracy using parallel k-anonymity (MapReduce model) algorithm defined in our prior work [1] and MapReduce based classification algorithm [2].

The rest of the paper is organized as pursues. Section 2 reviews related works. Section 3 presents proposed methodology. Section 4 provides evaluation metrics used. Section 5 presents experimental results. Section 6 concludes the paper while section 6 provides bearings for future work.

Revised Version Manuscript Received on 16 September, 2019.

* Correspondence Author

D. Radhika, Research Scholar, Computer Science Engineering, K L University, Guntur, AP, India. Email: radhikarajasekhar@yahoo.com

Dr D. Aruna Kumari, Professor, Department CSE, VJIT, Hyderabad, Telangana, India. Email: arunakumari@vjit.ac.in

II. RELATED WORKS

This section provides review of the prior work on huge information cleansing approaches. Bou-Harb *et al.* [3] proposed a model known as probabilistic darkness processing model for data sanitization. It also covers the concern on cyber situational awareness. It has inference probing in order to evaluate the solution and found to be effective in protecting privacy of big data. However, it has no provision for integrating the inference probes and the malware related issues. Waikar [4] on the other hand proposed a scalable and in-memory structure for huge information purification. They used Locality Sensitive Hashing (LSH) based solution which estimates risk prior to application of data sanitization. However, they found that risk based approach alone could not provide satisfactory solution.

Hemalatha and Elamparathi [5] discussed privacy issues in data mining practices and emphasized importance of preserving privacy. Like [3], on top of darkness professional mode, Aiswarya *et al.* [6] proposed a CTI based data sanitization model. It could provide a solution in the distributed environment. However, it just centered around the sanitization of misconfiguration data. Computation of sensitive and insensitive items is proposed in [7] for data sanitization. The solution is still to be implemented. Dasari and Rao [8] investigated different sanitization techniques that that could prevent inference attacks on social networks. Their research needs further enhancement to have useful solution. Soumya *et al.* [9] proposed a framework for big data sanitization. It was focusing on the controlled sanitization of big data in Hadoop environment. This work has limitations in experimental evaluation. Different privacy inference attacks and prevention measures are explored in [10].

Data sanitization and the process integrity were discussed in [11]. Their method monitors the sanitization process for integrity and ensures that it is done without security breaches. It is thus useful to have sanitization process correctly. However, no implementation details are provided. A more comprehensive review of big data privacy with various aspects is found in [12]. It clearly distinguishes between privacy and security. It throws light into different sanitization techniques. They opined that the research on privacy of big data is still open and needs different enhanced approaches. Natgunanathan *et al.* [13] investigated on protection of big data privacy. Their examination uncovered that the various methodologies provided solutions to protect big data. However, there are certain unresolved issues and challenges. Privacy is one of them which needs constant research and developments.

Rao *et al.* [14] considered privacy with respect to big data analytics. They could find certain threats in big data analytics. They include surveillance, disclosure of sensitive data, discrimination, abuse and personal embarrassment. They also discussed various privacy preserving methods like Multi-Dimensional Sensitivity Based Anonymization (MDSBA), techniques of cryptography, data distribution, randomization, T-closeness, L-diversity and k-anonymity. They proposed privacy preserving approach for data lake scenario. The

privacy challenges associated with big data are discussed in [16] while different solutions for the same are found in [17].

Different privacy related issues and solutions were found in [18], [19], [20] and [21]. From the review of literature, it is understood that the sanitization techniques were used to solve the problem of big data privacy. However, very few attempted the controlled anonymization and making well informed decisions while determining level of sanitization. In the second

III. PROPOSED FRAMEWORK

The framework discussed in this section was actually proposed in our prior work [2]. However, it is briefly discussed here prior to evaluating the misuse probability based big data sanitization approach. The big data provided to the framework is subjected to finding the attributes such as normal, sensitive and quasi identifiers. Sensitive identifiers reveal privacy details directly. Therefore, they are anonymized. Normal identifiers do not need privacy protection. However, quasi identifiers may be subjected to inference attacks. Therefore, it is essential to take care of the sanitization of data related to such identifiers. Novelty in our framework is to compute misuse probability score prior to deciding the degree of sterilization. in light of the score, it is possible to have certain degree of purification. This will help in avoiding undue sanitization that hinders utility of big data given for analytics. As one size does not fit all, the system recommends level of purification dependent on the given dataset and its attributes.

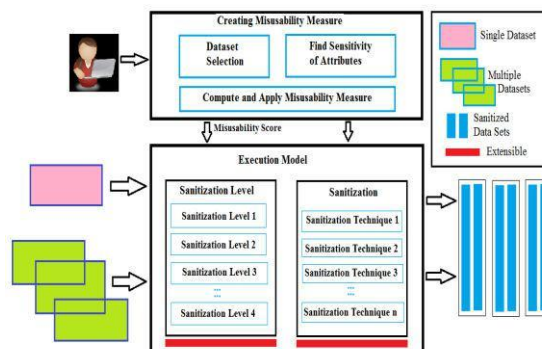


Fig 1: Overview of the framework

As shown in Figure 1, there are two important phases in the framework. The first is connected to finding misuse probability score which results in a value between 0.0 and 1.0. This measure can help in stage, in view of the misuse probability score, level of sanitization is applied appropriately instead of using sanitization techniques in general. This is where the framework plays crucial role which strikes balance between privacy of big data and the utility of the same when it is subjected to analytics. The latter is known as execution model which has two steps again. They include to determine level of sanitization and application of the sanitization appropriately to the given dataset.

3.1 Computing Misuse Probability Score

The process of computing misuse probability score is inspired by the work of Harel *et al.* [24] where extensive analysis is made on different kinds of attributes and the modalities to be used for data sanitization.

However, our prior works [1] and [2] and this paper focus on the enhancement and appropriate use of misuse probability score for big data in the context of MapReduce programming paradigm. Towards the empirical study, we defined the original k-anonymity algorithm using MapReduce model [1] which is used for evaluation of the proposed framework shown in Figure 1. Before moving to other details, the misuse probability score computation is discussed briefly here. However, more details of the same can be found in our previous work [2]. Computing misuse probability is a multi-step process. It includes computation of Raw Record Score (RRS) which reflects sensitivity of a tuple or record in given dataset which is tabular format, Record Distinguishing Factor (RDF) which is a measure to determine how far a quasi-identifier is capable of revealing identity or sensitive data and Final Record Score (FRS) which computes record score for entire dataset. Afterwards, the computation of misuse probability score is made. Eq. 1 is used to compute RRS.

After computing RRS, RDF is computed as in Eq. 2.

$$DF: \{ \text{quasi-identifiers} \} \in [0,1]$$

$$(2)$$

$$(3) \quad 0 \leq (R) \leq ()$$

The FRS is computed as in Eq. 3.

Finally misuse probability score is computed as in Eq. 4.

$$MS = \frac{1}{2} \times FRS = \frac{1}{2} \times ()$$

(4)

The score computed here is used by the proposed framework in order to have better possibilities in sanitization of big data while preserving its privacy and ensuring that utility is not lost unnecessary.

3.2 Misuse Probability Based Optimal Big Data Sanitization (MP-OBS)

Big data sanitization is made dependent on the misuse probability score in order to have appropriate level of sanitization. Anonymization may prompt loss of data utility. For this reason, it is given importance in this research to have proper degree of cleansing rather than applying privacy to big data using traditional approaches. As part of the framework shown in Figure 1, an algorithm is proposed to analyse and provide required decisions for anonymization or sanitization of big data. There is trade-off between privacy and utility and between misuse probability and level of sanitization. This trade-off is investigated for better results in big data privacy protection.

Algorithm: Misuse Probability Based Optimal Big Data Sanitization (MP-OBS)

Input: Dataset *D*

Output: Sanitized dataset *SD*

1. *level* = 0
2. *mps* = 0

3. *mps* = ComputeMPS(*D*) //Equations 1, 2, 3,4
4. IF *mps* >= 0.0 and <= 0.3 Then
5. *level* = 1;
6. For *i* in 1 to 50
7. *D*'[*i*] = Sanitize(*i*)
8. End For
9. Obtain Ground Truth with Data Analytics
10. *anonymity* = Determine correct anonymity value relating level 1
11. End If
12. IF *mps* >= 0.4 and <= 0.7 Then
13. *level* = 2
14. For *i* in 1 to 50
15. *D*'[*i*] = Sanitize(*i*)
16. End For
17. Obtain Ground Truth with Data Analytics
18. *anonymity* = Determine correct anonymity value relating level 2
19. End If
20. IF *mps* >= 0.8 and <= 1.0 Then
21. *level* = 3
22. For *i* in 1 to 50
23. *D*'[*i*] = Sanitize(*i*)
24. End For
25. Obtain Ground Truth with Data Analytics
26. *anonymity* = Determine correct anonymity value relating level 3
27. End If

Return *anonymity*

Algorithm 1: Misuse probability based optimal big data sanitization (MP-OBS)

As presented in Algorithm 1, the framework takes big data as input and the data is subjected to computation of misuse probability score. Based on the score, three conditions are used in order to know the degree of purification is needed. Each time, the dataset is subjected to different levels of anonymization to get ground truth from human experts. Finally, the algorithm is able to determine the level of sanitization based on the misuse probability score. The evaluation is carried out in order to have a standard way of dealing with the level of anonymization.

IV. EVALUATION METRICS

Different metrics are used for evaluation. They are known as Precision, Recall and Root Mean Square Error (RMSE). Precision is also known as Positive Predictive Value (PPV). Precision is the fraction of relevant results among obtained results. Similarly, recall refers to the fraction of relevant results over the total number of relevant results. RMSE on the other hand is the standard deviation associated with the prediction results. In other words, it is the measure that indicates how far the prediction errors are spread across the results. It is broadly used to find the difference between a model-predicted values and the actual values observed.



Table 1 shows confusion matrix from which precision and recall measures are derived.

	Ground Truth (correct prediction)	Ground Truth (wrong prediction)
Result of an algorithm (correct prediction)	TRUE Positive (TP)	FALSE Positive (FP)
Result of an algorithm (wrong prediction of intrusions)	FALSE Negative (FN)	TRUE Negative (TN)

Table 1: Confusion matrix

As presented in Table 1, there are different terms like True Positive (TP), False Positive (FP), False Negative (FN) and True Negative (TN). TP indicates that the algorithm has predicted something to be positive and in group truth also it is positive. TN means the algorithm correctly predicted all negative (incorrect) values. FP means the algorithm wrongly predicts as positive but the ground truth is negative. FN means that the algorithm wrongly predicts something as negative but in reality it is positive. Mathematical equations for Precision, Recall and RMSE are as in Eq.5, Eq. 6 and Eq. 7.

(5) Precision = True positive+False Negative

(6) Recall= (True positive+False Negative)/False Positive

$$(7) RMSE = \sqrt{\sum_{i=1}^n (P_i - O_i)^2}$$

RMSE is a standard measure where n represents number of observations, P denotes predicated value, O denotes observed value. These measures are used to evaluate the performance of given algorithm. In this paper, these measures are used for a classification algorithm that works with MapReduce framework. The algorithm is evaluated with the proposed framework where misuse probability based sanitization is explored.

V. EXPERIMENTAL RESULTS

Empirical study made with cloud environment such as Amazon Elastic Compute Cloud (compute engine which provides clusters), Amazon Elastic MapReduce (MapReduce environment provided by Amazon) and Amazon S3 (cloud based storage service). Two datasets from UCI machine learning repository [22] are used for the empirical study. Adult dataset and Census dataset are used for experiments. Classification algorithm written for MapReduce programming paradigm [23] is used to evaluate the proposed framework and the underlying algorithm.

5.1 Results with Adult Dataset

When big data is subjected to classification with anonymization as part of ground truth evaluation, the results are observed in terms of accuracy of classification, precision, recall and RMSE.

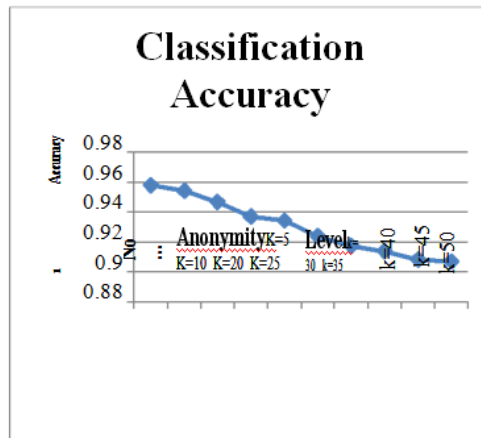


Fig 2: Classification accuracy vs. anonymity level

As presented in Figure 2, the level of anonymity is denoted in horizontal axis. The k value in the parallelized k-anonymity (out prior work in [1]) is taken from 5 to 50 for empirical study. The vertical axis shows classification accuracy. The results revealed that there in k value influence on the accuracy of classifier [23]. When k value is increased, the classification accuracy is reduced due to the reduced utility of dataset as a result of anonymization.

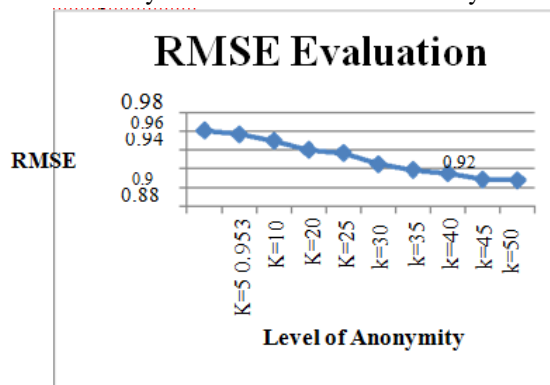


Fig 3: RMSE evaluation

As shown in Figure 3, the degree of obscurity is taken from 5 to 50 in horizontal axis and vertical axis shows RMSE which reveals distribution of error. It is understood that the k value used in parallelized k-anonymity algorithm (our prior work in [1]) which is used for anonymization, has its impact on the RMSE. RMSE is decreased when level of anonymity is increased.

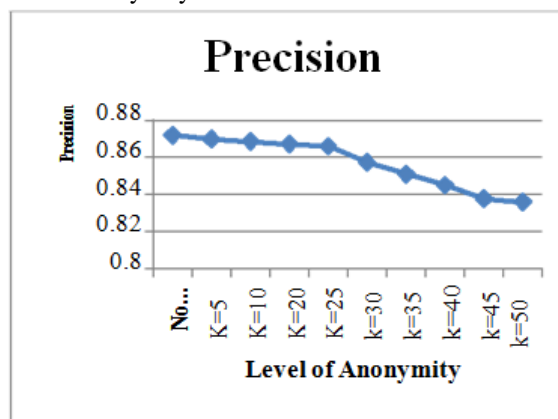


Fig 4: Level of anonymity vs. precision

As presented in Figure 4, the precision analysis is made. The level of anonymity is provided in horizontal axis while the vertical axis presents the precision value. As discussed in the preceding section precision is one of the measures to know the performance of an algorithm. The precision exhibited by the classification algorithm [23] is decreased when there is increase in the level of anonymity.

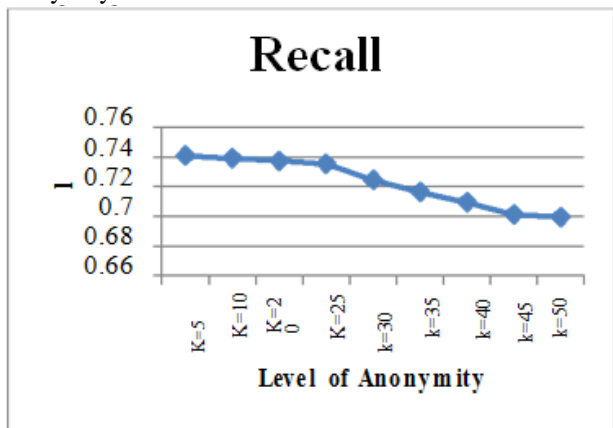


Fig 5: Recall value vs. level of anonymity

As presented in Figure 5, it is evident that the k value for parallel anonymization technique (our prior work in [1]) is presented in horizontal axis while the vertical axis shows the recall value. The results revealed that the recall value is decreased when level of anonymity is increased. The recall reflects the performance of the classification algorithm (MapReduce version from [23]) with different level of sanitization.

5.2 Results with Census Dataset

This sub section provides results of accuracy, precision, recall and RMSE of classification algorithm (MapReduce version [23]) with Census dataset.

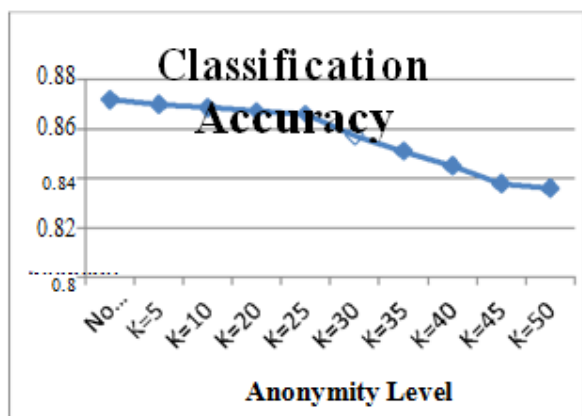


Fig 6: Anonymity level vs. classification accuracy

As presented in Figure 6, the level of anonymity is denoted in horizontal axis. The k value in the parallelized k-anonymity (our prior work in [1]) is taken from 5 to 50 for empirical study. The vertical axis shows classification accuracy. The results revealed that there is k value influence on the accuracy of classifier [23]. When k value is increased, the classification accuracy is reduced due to the reduced utility of dataset as a result of anonymization.

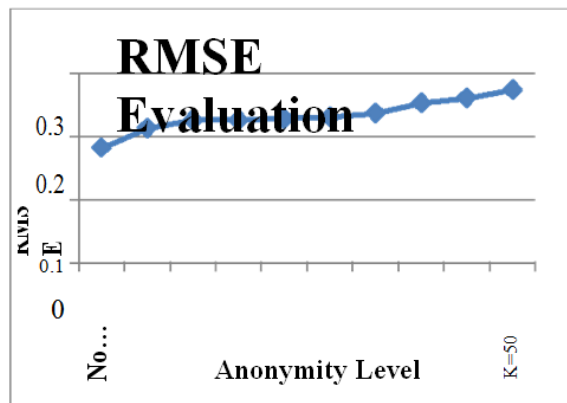


Fig 7: RMSE evaluation

As shown in Figure 7, the level of anonymity is taken from 5 to 50 in horizontal axis and vertical axis shows RMSE which reveals distribution of error. It is understood that the k value used in parallelized k-anonymity algorithm (our prior work in [1]) which is used for anonymization, has its impact on the RMSE. RMSE is slightly increased when level of anonymity is increased.

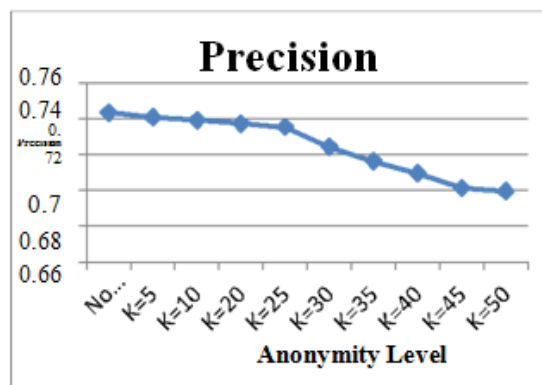


Fig 8: Level of anonymity vs. precision

As presented in Figure 8, the precision analysis is made. The level of anonymity is provided in horizontal axis while the vertical axis presents the precision value. As discussed in the preceding section precision is one of the measures to know the performance of an algorithm. The precision exhibited by the classification algorithm [23] is decreased when there is increment in the degree of secrecy.

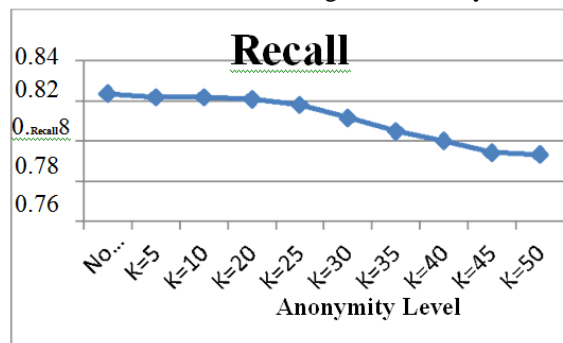


Fig 9: Recall value vs. level of anonymity

As presented in Figure 9, it is evident that the k value for parallel anonymization technique (our prior work in [1]) is presented in horizontal axis while the vertical axis shows the recall value.

The results revealed that the recall value is decreased when level of anonymity is increased. The recall reflects the performance of the classification algorithm (MapReduce version from [23]) with different level of sanitization.

VI. CONCLUSION AND FUTURE WORK

Misuse of sensitive data has been a problem with insider attacks or privacy attacks made by external adversaries. The severity of the problem is increased with big data being outsourced to public cloud which is considered an untrusted environment. Existing data sanitization or anonymization techniques are useful to protect privacy of big data. However, this may lead to losing utility of data. Therefore, striking balance between sanitization and utility is essential. Towards this end, we planned a system that has provision for finding misusability probability of given dataset. Based on the probability of misuse, the framework determines the level of sanitization which will preserve utility of the data besides protecting its privacy. The exploration in this paper centers around the analysis of trade-offs between misuse probability and level of sanitization. Towards this end, an algorithm is proposed. The solution is implemented with real cloud environment which includes Amazon EC3, S3 and EMR. The assessment could give helpful bits of knowledge on the exchange offs. The empirical results revealed that the proposed framework is useful in the real world scenarios associated with big data publishing and data analytics on big data. The framework is evaluated with tabular data or structured data. It does not support unstructured data. It is an important drawback as the big data includes structured, semi-structured and unstructured data. This will be examined in our future work.

REFERENCES

- [1] D. Radhika and D. Aruna Kumari. (2016). A Framework for Exploring Algorithms for Big Data Mining. *IJST*. 9 (17), p1-7.
- [2] D. Radhika and D. Aruna Kumari. (2018). Misusability Measure Based Sanitization of Big Data for Privacy Preserving MapReduce Programming. *IJECE*. 8 (6), p4524-4532.
- [3] Bou-Harb, E., Husak, M., Debbabi, M., & Assi, C. (2017). *Big Data Sanitization and Cyber Situational Awareness: A Network Telescope Perspective*. *IEEE Transactions on Big Data*, 1-17.
- [4] Kanchan Prakash Waikar. (2017). BIG DATA SANITIZATION USING SCALABLE IN-MEMORY FRAMEWORKS. *MASTER OF SCIENCE IN COMPUTER SCIENCE*, p1-110.
- [5] R. Hemalatha and M. Elamparithi. (2015). Privacy Preserving Data Mining Using Sanitizing Algorithm. *IJCS and Information Technologies*. 6 (5), p4174-4179.
- [6] S. Aiswarya, S. Usharani, K. Dhanalakshmi and M. Roberts Masillamani. (2018). A CTI BASED BIG DATA SANITIZATION USING DARKNET PREPROCESSING MODEL. *International Journal of Pure and Applied Mathematics*. 119 (14), p1095-1100.
- [7] S. Devika M. Sc and Dr. D. Devakumari. (2016). Data Sanitization for Preserving Data Privacy. *IJAR in Computer Engineering & Technology*. 5 (11), p1-9.
- [8] JAYASREE DASARI and K.R. KOTESWARA RAO. (2014). Sanitization Techniques for Protecting Social Networks from Inference Attacks. *International Journal of Computer Science and Mobile Computing*. 3 (12), p 236-244.
- [9] Y. SOWMYA, Dr. M NAGARATNA and Dr.C SHOBA BINDU. (2005). M-SANIT: A FRAMEWORK FOR EFFECTIVE BIG DATA SANITIZATION USING MAP REDUCE PROGRAMMING IN HADOOP. *Journal of Theoretical and Applied Information Technology*. 96 (6), p1-10.
- [10] Raymond Heatherly, Murat Kantarcioglu, and BhavaniThuraisingham. (2013). Preventing Private Information

- Inference Attacks on Social Networks. *TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING*. 25 (8), p1-14.
- [11] Bujwala and P Raja Sekhar Reddy. (2016). An Effective Mechanism for Integrity of Data Sanitization Process in the Cloud. *European Journal of Advances in Engineering and Technology*. 3 (8), p82-84.
- [12] Jain, P., Gyanchandani, M., & Khare, N. (2016). *Big data privacy: a technological perspective and review*. *Journal of Big Data*, 3(1). P1-25.
- [13] Mehmood, A., Natgunanathan, I., Xiang, Y., Hua, G., & Guo, S. (2016). *Protection of Big Data Privacy*. *IEEE Access*, 4, 1821-1834.
- [14] Ram Mohan Rao, P., Murali Krishna, S., & Siva Kumar, A. P. (2018). *Privacy preservation techniques in big data analytics: a survey*. *Journal of Big Data*, 5(1). P1-12.
- [15] Nitin Kumar Agrawal and Aprna Tripathi. (2015). Big Data Privacy Challenges and Techniques. *International Journal of Computer Applications*, p1-5.
- [16] Altman, M., Wood, A., O'Brien, D. R., & Gasser, U. (2018). *Practical approaches to big data privacy over time*. *International Data Privacy Law*, 8(1), 29-51.
- [17] Dr.C. Nalini and Dr.A.R. Arunachalam. (2017). A STUDY ON PRIVACY PRESERVING TECHNIQUES IN BIG DATA ANALYTICS. *IJPAM*. 116 (10), p281-286.
- [18] Jain, P., Gyanchandani, M., & Khare, N. (2016). *Big data privacy: a technological perspective and review*. *Journal of Big Data*, 3(1).
- [19] Dr. Puneet Goswami and Ms. Suman Madan. (2017). A Survey on Big Data & Privacy Preserving Publishing Techniques. *Advances in Computational Sciences and Technology*. 10 (3), p 395-408.
- [20] John P. Holdren and Eric S. Lander. (2014). REPORT TO THE PRESIDENT BIG DATA AND PRIVACY: A TECHNOLOGICAL PERSPECTIVE. *President's Council of Advisors on Science and Technology*, p1-76.
- [21] Bao, R., Chen, Z., & Obaidat, M. S. (2018). *Challenges and techniques in Big data security and privacy: A review*. *Security and Privacy*, p1-8.
- [22] UCI (2017). UCI Machine Learning Repository. Available online at: <https://archive.ics.uci.edu/ml/index.php>. [accessed on: 20 December 2018]
- [23] William Cohen (2017). Naïve Bayes and Map-Reduce. Retrieved from <https://pdfs.semanticscholar.org/a397/eb310921897ef8a140668b623de618da7606.pdf>. Accessed on 20 January 2019.
- [24] Harel, A., Shabtai, A., Rokach, L., and Elovici, Y. (2012). M-Score: A Misusability Weight Measure. *IEEE Transactions on Dependable and Secure Computing*, 9 (3), p414-428.

AUTHORS PROFILE



Mrs D Radhika is Asst Professor at Department of CSE, Stanley College of Engineering & Technology for Women, Hyderabad, Telanagana, India. Data mining, Big Data Analytics, Internet of Things, Privacy of Big Data are the interested areas of research. She published 7 research articles in International Journals. She also presented and published research papers in National and International Conferences. She is currently pursuing research on "Privacy of Big Data".



Dr Aruna Kumari Professor at Department of CSE,VJIT, Hyderabad, Telangana India. She is Fellow of CSI(FCSI), and Fellow of IEEE (FIEEE). She is DST Young Scientist Awardee (Govt. of India). She is Ph.D Supervisor at Department of CSE KL-University, Guntur, Andhra Pradesh. Data mining, Cloud Computing, Privacy & Security of Big Data, Big Data Analytics are her interested areas of research. She has more than 70 research articles in International Journals and Conferences. She is Currently working on research project "Design and Development of Effective Privacy Preserving Data Mining for Cardiac Cancer and Diabetic Health Care" funded by DST-SERB. She has organized several National and International Conferences and Workshops.