

Privacy Preserving Data Mining Technique - Issues and Challenges

Mayur Rathi

Department of Computer Science & Engineering
Shri Vaishnav Vidyapeeth Vishwavidyala, Indore, India
e-mail: mayurrathi.31dec@gmail.com

Anand Rajavat

Department of Computer Science & Engineering
Shri Vaishnav Vidyapeeth Vishwavidyala, Indore, India
e-mail: anandrajavat@yahoo.co.in

Abstract— The technology improves the ability of data storage and their processing. Therefore a significant growth on different business domains is observed. The business intelligence requires data and consumer habits but the requirement of data is never fulfilled by a single data source. In this context required to combine the different source of information for making smart decisions or for finding the conclusions. But the limitation is that nobody wants to disclose their client's private and sensitive information to others. Therefore we need a privacy preserving environment to handle the data, it's privacy and conclusion. In this presented paper privacy preserving data mining is the key area of investigation using the available literature. Additionally by concluding the literature need to establish the future prospects for the work.

Keywords— *Privacy preserving data mining; multi party data source; objective summarization; model improvements; techniques and tools.*

I. INTRODUCTION

Data mining is an application of analyzing data, for extraction of meaningful patterns. That helps to understand and distinguish data variations in different applications. The use of data mining is also performed for making decision and predicting the future aspect of growth or downfall. Classically the utilization of data mining approaches is found in various centralized data bases. Data is stored in one place and the data mining algorithms and techniques are applied to perform the application centric task [1].

There is another scenario where data are distributed in different data sources. Additionally the data mining algorithms processes data for recovering the desired facts form data without spoiling privacy and confidentiality of data. Therefore this type of data analysis technique is known as the Privacy Preserving Data Mining (PPDM) [2].

Principally, the privacy preserving data modeling is utilized where different parties are providing data and associated for concluding a common solution. Even though, no one wants to expose their data to other

party. In this context the available data may be in form of [3]:

- Vertical partitioned data
- Horizontal partitioned data

In addition of that it may be possible all parties can faith on the common party to analyze or not therefore the environment can also affect the processing of data.

- Semi trusted environment
- Fully trusted environment

In addition of that there are other issues and challenges which can affect the process of privacy preserving data mining and their techniques. Therefore in this paper the detailed study of privacy preserving data modeling is conducted. Additionally on the basis of collected literature the future objectives are established.

II. BACKGROUND

This section incorporates the relevant information and key terms which are used during the proposed privacy preserving data mining domain.

A. Data Mining

Data mining is an art of handling data in order to discover the application oriented patterns and decision discovery. Therefore, it is a rich domain of applications and possibilities. A number of applications now in these days utilizing the services of data mining for improving their production, business directions, market analysis and more [4]. There are two main classes of data mining algorithms supervised and unsupervised. In supervised learning some predefined samples are required to learn on the patterns. After successful learning the supervised algorithms are able to recognize the similar patterns. On the other hand the unsupervised algorithms not required to train with the predefined samples, these algorithms are directly applicable to the data for discovering patterns.

B. Privacy Preserving

Now in these days every business domain gathering ends client's data. Among the gathered data some of part of data is much confidential, sensitive and private. Additionally, discloser of such kind of data to any other party can harm the end client privacy socially and financially. Therefore, during the process of data pattern discovery, it is required to preserve the private and confidential information for experimentation and other analytical task in any business or research domain [5]. In other terms keep the sensitive data confidential from the unauthorized persons and environment is known as the privacy preserving.

C. Data Discloser

As discussed in privacy preserving scenario some of the data is sensitive, private and confidential such kind of data distribution can harm the end client. Therefore before discloser of data for public or experimental use it is required to transform the data's actual figure, encrypt the data or encode data to non understandable format is necessary. But data discloser can be done either intentionally or by mistake. Therefore, handling of such kind of data is much essential [6]. Therefore, in some of the cases the data sensitivity scanning and their normalization is essential part of security quality of service. After that process the data is dis-closeable for public media.

D. Dimensionality Reduction

The privacy preserving Data mining techniques requires data from different parties and data sources.

Additionally, the data is clubbed from different sources in terms of attributes and the number of instances. Therefore, it has a significant dimension of data. In addition of that as the number of parties increases the dimensions of data can be increased. The processing of higher dimensional data requires the significant amount of processing cost i.e. time and space complexity. Therefore it is required to identify the essential features from the bulk set of data to minimize the data processing cost and improving the performance of classification [7].

E. Impact of Noise on Data

We try to explain an example for describing this context. A person attending phone call in market or in high traffic area, it is much complicated to understand all the communication for that person. On the other hand the same person attending phone in their house minimizes the effort to understand the phone call communication. In this context noise can affect the person's perception. In the same way when a machine learning algorithm tries to learn on noisy data it is directly impact on classifier's or learning algorithm's performance. Therefore, it is required to identify the high informative features among entire information available on datasets [8].

F. Classification

Basically, the classification is a supervised learning technique in data mining and machine learning. In such kind of pattern learning techniques the two sets of data required first set of data is termed as training samples or patterns. These training patterns are predefined observations with their specific outcomes. The learning algorithms or classifiers are learnt on these samples and prepare the model. To recognize the similar pattern on which the learning algorithm trained on [9]. The classification algorithms can be transparent in nature or it may be opaque. The transparent algorithms develop rules and by invoking these rules the decisions are calculated. On the other hand in opaque algorithms the weights or other factors are calculated for recognition for unknown patterns.

G. Decision Making

Decision making process can be supervised and unsupervised in nature. A significant amount of literature is available where the unsupervised learning techniques are also used for making decisions using the available data. But in our context the decision making is a conclusion based on the input set of data instances. Therefore, employment of both kinds of

algorithms can be possible for making decision and prediction for a given unknown pattern of data instance [10]. In a number of applications, the visualization techniques are used for decision making i.e. heat map, decision trees. In addition of sometimes data mining algorithms are applied. When a journal decision required the visualization, techniques are effective way to preserve the time. On the other hand, the data mining algorithms can be applied when the precise outcomes are expected.

H. Data Collaboration

In business intelligence domain sometimes it is not feasible to take decision by using incomplete set of information. Therefore, it is required to delegate the additional information from other data source or party to complete the set of knowledge. For finding the essential patterns, to take big decisions in business domain required to complete analysis of data. The complete set of knowledge belongs to similar industries help to make smart decisions. Therefore, the different business owners agreed to aggregate their data for making common decisions. The delegation of data or information also needs to be secured for preserving the privacy of data [11].

III. LITERATURE SURVEY

The privacy preserving data mining is one of the promising areas in data science. In order to improve the business and other business intelligence task there are a number of scenarios where we need to handle data sources from different data sources. Additionally, it is also required to combine and explore data to find common fruitful outcomes. The Privacy Preserving Data Modeling is the only method which solves the issue of data privacy and data utility in different situations.

For instance, consider a tour and travel industry where the hotels, restaurants, cabs, and other transportation mediums are combined each other. Additionally, the information about their clients is needed to be preserve by the service providers. In this context if they want to combine their data and want the combined conclusion for any research purpose then security and privacy on data mining is necessary [12].

Consider separate medical institutions that wish to conduct a joint research while preserving the privacy of their patients. One way to view this is to imagine a trusted third party-- everyone gives their input to the

trusted party, who performs the computation and sends the results to the participants. However, this is exactly what we don't want to do, for example, hospitals are not allowed to hand their raw data out, security agencies cannot afford the risk, and governments risk citizen outcry if they do. Thus, the question is how to compute the results without having a trusted party, and in a way that reveals nothing but the final results of the data mining computation. Secure Multiparty Computation enables this without the trusted third party [13].

Chen et al. [14] conducted a review on data intensive techniques and their applications. According to this survey when the data is available in higher dimension then for classifiers performance becomes low in this context the technique employed for Dimensionality reduction can be unsupervised [15] and supervised [16] in nature. Sunitha et al. [17] perform experiments on noisy data and outlier data according to their study both the issues impact on classifiers performance. In this context Xiong et al. [18] suggests the technique of noise removal in data. In further Zhang et al. [19] shows the issues of privacy in multiparty data collaboration. Therefore it is required to deal with multiparty data with sensitive information for finding a secure computation and analysis of the data. Therefore the proposed work is intended to address the key issues and challenges involve in privacy preserving data mining efficiency and the performance of mining.

IV. RESEARCH GAP ANALYSIS FOR EXISTING METHODS

This section provides the extracted essential research gap for concluding the proposed work.

Table 1: Research Gap analysis of existing methods

Author, Publication, Year	Research Gap
Kung et al. [20], ACM Transactions on Embedded Computing Systems, July 2017	The given technique only deals with the dimensionality issues and security concern with the data in specific scenario.
Zhang et al. [21], IEEE Transactions On Computers, May 2016	The method is not providing the solution for data publishing to any third party with enhanced utility and assurity
Nissim et al. [22], Information Science, Elsevier, 2010	Over flexibility, over anonymity, over computational cost & also data projection based

	approach is followed with is not suitable for all the kind of data in real world
Poovammal et al. [23], Journal of Computer Science (JCS), 2009	The given method only works for the numerical attributes and for categorical data it is not functioning
Li et al. [24], IEEE Transactions on Knowledge and Data Engineering, March 2012	works only for the high-dimensional data not considering the PPDM scenario
Xiong et al. [18], IEEE Transactions on Knowledge and Data Engineering, March 2006	All the methods perform well in specific scenario. Author suggested working on combination of two or three methods for generalization.
Bouzas et al. [16], IEEE Transactions on Neural Networks and Learning Systems, May 2015	Proposed algorithm is not suitable for very high dimensional data
Fletcher et al. [25], International Journal of Computer Theory and Engineering (IJCTE), February 2015	Authors are not suggesting a specific technique to work and enhance the quality of information.
Hua et al. [26], IEEE Transactions on Information Forensics and Security, October 2016	Author are not concerned to deal with the data utility in place of data query able
Li et al. [27], IEEE Transactions on Information Forensics and Security, 2016	Proposed scheme increases the complexity of system in terms of resources
Li et al. [28], IEEE Transactions on Knowledge and Data Engineering September 2012	Authors are not considering the rust of data publishing scenario.
Arribas et al. [29], Information Processing & Management, 2012	Authors are not considering the issue of data aggregation and mining security

V. ISSUES AND CHALLENGES

Data mining techniques and applications are growing continuously. Additionally, its acceptability is also improving day by day in various real-world applications for decision making and another context. But sometimes the available data in a single source is not much sufficient to deal with the issues, therefore, it is required to gather information from other sources, to develop a complete dataset for the decision making

the process. Thus the third party data contributor is worried about the data and consumers' privacy, sensitivity and confidentiality due to the risk of information discloser. In this context using the existing literature some key issues are concluded as:

1. **Network Data security issues:** in the PPDM environment, multiple parties are agreed to combine their data for finding the common data model and the decision ability. Thus the data is communicated between parties and the centralized server. The communication between servers and parties is enabled using public networks. Additionally, public networks are not secure thus there is a probability for network security issues during the communication of sensitive and confidential data using network attacks.
2. **Data discloser and data leakage issue:** although the data in PPDM is modified, transformed and/or sometimes encrypted to preserve the privacy and data owners' security. But due to the mistake of algorithm or computational losses, a transformed or modified data can provide the values of a data owner accidentally.
3. **Data noise and outcomes:** in PPDM to preserve the privacy and confidentiality of data, noise is introduced. But, the noise in data can affect the performance of the learning algorithm. Additionally, it can also impact on outcomes of the learning algorithm. Therefore, it is required to introduce the noise in data in a controlled manner.
4. **Dimensionality:** the PPDM collaborate the data of a number of parties. The collaboration of data can be possible in a horizontal or vertical manner. In both, the conditions parties involve their data by which either the number of instances is increases or the number of attributes increases. Therefore both the manner of data aggregation increases the dimensionality of data.
5. **Data utility:** in data mining applications data is processes in their absolute format and the outcome of data mining is used to get benefits for the application. But in PPDM the data is used after modification in the actual values.

Therefore the outcome of the data mining can also be noisy and the issue of data utility rises.

6. **Data publishing:** in PPDM the data discloser and publishing are also performed. Therefore, to use data in other application dataset is modified with some noise for security and privacy purposes. The modification of data can impact the application's performance, the application conclusions, and decision-making ability.
7. **Mining environment:** that is an essential question in PPDM. How the data is secure, safe and confidential in third party storage. The data parties are always worried about the trust level of data processing and storage. What happened when the data processing server is compromised with the attacker? This serves the purpose of intention.

In order to achieve an effective PPDM model, we have to consider different aspects of PPDM such as working with the Efficient Data Dimensionality Reduction Techniques to preserves the computational resources. Therefore, investigation of different Dimensionality Reduction Techniques required. On the other hand, the learning algorithm's performance depends on the quality of data, therefore, need to Measure the Impact of Noise over Classifier Performance. Moreover, there are some approaches available that are preserving the privacy of data using noise. Thus, a suitable and controlled noise composition required to maintain the utility of data.

VI. PROPOSAL FOR PRIVACY PRESERVING DATA MODEL

In previous section the review of existing work is reported and based on the literature collection a probable model is presented in this section.

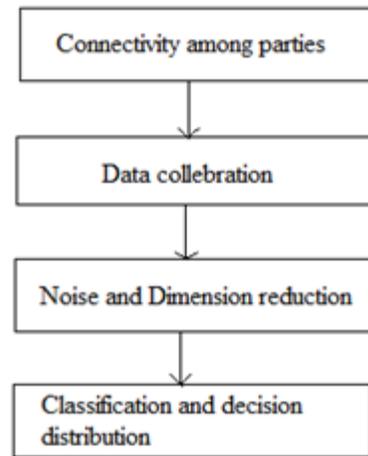


Fig. 1. Praposed Data Model

A. *Connectivity among Parties*

The proposed data model is aimed to satisfy the requirements of privacy preserving in data mining. Therefore, it may be possible that there are N number of parties exist, which want to combine their data for finding effective decisions. Therefore, this phase of the system is responsible for the following two expectations:

- Connectivity of available parties with the semi trusted server
- Enable security and privacy at the client end for regulating the part of data which are need to be disclosed and essential for decision making. Therefore this phase is also enable the cryptography, data encoding, data transformation or mapping to securing the sensitive part of data. This method will solve the purpose as per pridictions.

B. *Data Collaboration*

The combination of data is required in this scenario, but the synchronization of data is also a significant need of the system. The un-synchronized data can misguide the learning algorithm and produces undesired results. Therefore, this phase of system is dedicated to verifying the data instances and their relevance class labels.

C. *Noise and Dimensionality Reduction*

After synchronization of data it is required to verify the quality of data and availability of information of the present attributes to reduce the impact of noise and cost of processing for large dimensional data. Therefore in this phase the noise

measurement technique or feature selection technique is required to optimize the quality of data.

D. Classification and decision distribution

That is the final stage of the proposed system. This phase is responsible for mining that data, which is obtained from the previous phase. Using this data, the required application-oriented decisions are developed. Additionally, the data decisions are distributed to all the participating parties who are providing the data for experimentation and decision mining. During the distribution of final decisions and experimental data set disclosure it is ensured that the sensitive and private attributes from the data set keep preserved. The attributes of any data are the issues related to originality and will help in verification.

VII. CONCLUSION AND FUTURE WORK

Privacy preservation in data mining has become more important in recent years, because of the increasing ability to store personal data about users, and the increasing sophistication of data mining algorithms to leverage this information. If information management or possession is distributed, the disclosure of personal data becomes a problem. The privacy preserving data mining leads a number of issues and research challenges, among them accurately data processing is a key research issues. During such kind of data modeling research and their publishing the following situations can affect the modeling of data:

- Data nature and dimensions
- Data integration processes and affect of modeling
- Data publishing and quality of outcomes

Privacy Preserving Data Model is a multiparty data mining environment where from multiple data sources data is combined to achieve a common objective. So the aim of proposed work is to preserve the expected level of privacy with minimum information loss, with low error rate and improvement in accuracy in multiparty data sets. In order to achieve the mentioned aim, the future objectives of the proposed work are established as:

- 1) The Efficient Data Dimensionality Reduction Techniques preserves the computational resources therefore need to investigate different Dimensionality Reduction Techniques for proposed system.

- 2) The learning algorithms performance depends on the quality of data which is used for training therefore need to Measure the Impact of Noise over Classifier Performance.
- 3) There are some approaches available that are preserving privacy of data using noise thus what is suitable and effective noise composition required to maintain the utility of data is need to be investigate.
- 4) The aim of PPDM techniques is to combine and mine the data. Additionally discover the patterns which are useful for all the participating parties. Thus need to Design and Develop the Efficient Common Data Publishing Technique.

REFERENCES

- [1] A. Mukhopadhyay, U. Maulik, S. Bandyopadhyay, C. A. Coello, "A survey of Multi-objective Evolutionary Algorithms for Data Mining: Part I", *IEEE Trans. on Evolutionary Computation.* 18(1), pp. 20-35, 2014.
- [2] Y. A. A. S. Aldeen, M. Salleh, M. A. Razzaque, "A comprehensive review on Privacy Preserving Data Mining" *Springer Plus* 4:694, pp. 1-36, 2015.
- [3] J. Danasana, R. Kumar, D. Dey, "Mining Association Rules for Horizontally Partitioned Databases using CK Secure Sum Technique", *Int. J. of Dist. and Parallel Sys.* 3(6), pp. 149-157, 2012.
- [4] J. Ponce, A. Karahoca, "Data Mining and Knowledge Discovery in real life applications", In-Teh, Croatia, 2009.
- [5] R. Mendes, J. P. Vilela, "Privacy-Preserving Data Mining: Methods, Metrics, and Applications", *IEEE Access.* 5, pp. 10562-10582, 2016.
- [6] L. Xu, C. Jiang, J. Wang, J. Yuan, Y. Ren, "Information security in Big Data", privacy and data mining. *IEEE Access.* 2, pp. 1149-1176, 2014.
- [7] K. K. Vasan, B. Surendiran, "Dimensionality Reduction using Principal Component Analysis for Network Intrusion Detection", *Perspectives in Science.* 8, pp. 510-512, 2016.
- [8] D. F. Nettleton, A. Orriols-Puig, A. Fornells, "A study of the effect of different types of Noise on the precision of Supervised Learning Techniques", *Artif Intell Rev.* 33(4), pp. 275-306, 2010.
- [9] P. Anitha, G. Krithka, M. D. Choudhry, "Machine Learning Techniques for learning features of any kind of data: A Case Study", *Int. J. of Adv. Research in Comp. Engg. & Tech.* 3(12), pp. 4324-4331, 2014.
- [10] B. Chitradevi, N. Thinaharan, "Role of decision making in Data Mining Systems", *Int. J. of Trend in Research and Dev.,* 2(5), pp. 122-125, 2015.
- [11] S. K. Swamy, S. H. Manjula, K. R. Venugopal, S. S. Iyengar, L. M. Patnaik, "Association rule sharing Model for Privacy Preservation and Collaborative Data Mining efficiency", In: *RAECS 2014 Proceedings of 2014 Recent Advances in Engineering and Computer Sciences*, pp. 1-6, IEEE, Chandigarh, India, 2014.
- [12] A. Basiri, P. Amirian, A. Winstanley, T. Moore, "Making Tourist Guidance Systems more Intelligent, Adaptive and Personalised using Crowd Sourced Movement Data", *J*

- Ambient Intell Human Comput, 2017, 9:413. <https://doi.org/10.1007/s12652-017-0550-0>.
- [13] G. Kou, Y. Peng, Y. Shi, Z. Chen, "Privacy-Preserving Data Mining of Medical Data using Data Separation-Based Techniques", *Data Science Journal*, 6, Suppl., pp. 429-434, 2007.
- [14] C. L. P. Chen, C. Y. Zhang, "Data-Intensive Applications, Challenges, Techniques and Technologies: A survey on Big Data", *Information Sciences* 275, pp. 314-347, Elsevier, 2014, <https://doi.org/10.1016/j.ins.2014.01.015>.
- [15] C. Xu, D. Tao, C. Xu, Y. Rui, "Large-Margin Weakly Supervised Dimensionality Reduction" In: *Proceedings of the 31st International Conference on Machine Learning*, pp. 865-873. JMLR: W&CP Beijing, China, 2014.
- [16] D. Bouzas, N. Arvanitopoulos, A. Tefas, "Graph Embedded Nonparametric Mutual Information for Supervised Dimensionality Reduction", *IEEE Trans. on Neural Networks and Learning Systems*, 26(5), pp. 957-963, 2015.
- [17] L. Sunitha, M. B. Raju, B. S. Srinivasa, "A Comparative Study between Noisy Data and Outlier Data in Data Mining" *Intl. J. of Current Engg. and Tech.*, 3(2), pp. 575-577, 2013.
- [18] H. Xiong, G. Pandey, M. Steinbach, V. Kumar, "Enhancing Data analysis with Noise removal" *IEEE Trans. on Knowledge and Data Engineering*, 18(3), pp. 304-319, 2006.
- [19] W. Zhang, Y. Lin, S. Xiao, J. Wu, S. Zhou, "Privacy Preserving Ranked Multi-Keyword search for multiple data owners in Cloud Computing", *IEEE Trans. on Computers* *Journal of Latex Class Files*, 6(1), pp. 1-14, 2015.
- [20] S Y Kung, T Chanyaswad, J M Chang, P Wu, "Collaborative PCA/DCA Learning Methods for Compressive Privacy", *ACM Transactions on Embedded Computing Systems*, Vol. 16, Issue 3, Article 76, pp. 77-117, July 2017.
- [21] Q Zhang, L T. Yang, and Z Chen, "Privacy Preserving Deep Computation Model on Cloud for Big Data Feature Learning", *IEEE Transactions On Computers*, Vol. 65, No. 5, pp. 1351-1362, May 2016.
- [22] M Nissim, L Rokach, and O Maimon. "Privacy-Preserving Data Mining: A Feature Set Partitioning Approach", *Information Sciences* 180, Elsevier, No. 14, pp. 2696-2720, 2010.
- [23] E. Poovammal and M. Ponnaivaikko, "Utility Independent Privacy Preserving Data Mining on Vertically Partitioned Data", *Journal of Computer Science*, Volume 5, Issue 9, pp. 666-673, 2009.
- [24] T Li, N Li, J Zhang and I Molloy, "Slicing: A New Approach for Privacy Preserving Data Publishing", *IEEE Transactions on Knowledge and Data Engineering*, Volume: 24, Issue 3, pp. 561-574, March 2012.
- [25] S Fletcher and M Z Islam, "Measuring Information Quality for Privacy Preserving Data Mining", *International Journal of Computer Theory and Engineering*, Vol. 7, No. 1, pp. 21-28, February 2015.
- [26] J Hua, A Tang, Y Fang, Z Shen, and S Zhong, "Privacy-Preserving Utility Verification of the Data Published by Non-interactive Differentially Private Mechanisms", *IEEE Transactions on Information Forensics and Security*, Doi 10.1109/TIFS.2016.2532839, Volume 11, Issue 10, pp. 2298-2311, Oct 2016.
- [27] L Li, R Lu, K K R Choo, A Datta, and J Shao, "Privacy-Preserving Outsourced Association Rule Mining on Vertically Partitioned Databases", *IEEE Transactions on Information Forensics and Security*, DOI 10.1109/TIFS.2016.2561241, pp. 1847-1861, 2016.
- [28] Y Li, M Chen, Q Li and W Zhang, "Enabling Multi-level Trust in Privacy Preserving Data Mining", *IEEE Transactions on Knowledge and Data Engineering*, Volume 24, Issue 90, pp. 1598-1612, Sep 2012.
- [29] N Arribas, G V Torra, A Erola, and J C Roca. "User k-Anonymity for Privacy Preserving Data Mining of Query Logs", *Information Processing & Management*, Volume 48, No. 3, pp. 476-487, 2012.