# Performance analysis of perturbation-based privacy preserving techniques: an experimental perspective

**Ritu Ratra, Preeti Gulia, Nasib Singh Gill**

Department of Computer Science and Applications, Maharshi Dayanand University, Rohtak, India

## Article Info

## ABSTRACT

Nowadays, enormous amounts of data are produced every second. These data also contain private information from sources including media platforms, the banking sector, finance, healthcare, and criminal histories. Data mining is a method for looking through and analyzing massive volumes of data to find usable information. Preserving personal data during data mining has become difficult, thus privacy-preserving data mining (PPDM) is used to do so. Data perturbation is one of the several tactics used by the PPDM data privacy protection mechanism. In perturbation, datasets are perturbed in order to preserve personal information. Both data accuracy and data privacy are addressed by it. This paper will explore and compare several hybrid perturbation strategies that may be used to protect data privacy. For this, two perturbation-based techniques named improved random projection perturbation (IRPP) and enhanced principal component analysis-based technique (EPCAT) were used. These methods are employed to assess the precision, run time, and accuracy of the experimental results. This paper provides the impacts of perturbation-based privacy preserving techniques. It is observed that hybrid approaches are more efficient than the traditional approach.

*Corresponding Author:*

Ritu Ratra
Department of Computer Science and Applications, Maharshi Dayanand University
Rohtak, Haryana, India
Email: ritu.rs.dcsa@mdurohtak.ac.in

## 1. INTRODUCTION

Mining datasets spread across numerous organizations without disclosing further personal information has recently gained importance. In many businesses, protecting data privacy is currently a big concern. In present scenario, people are worried that their private information may be exposed and used for improper purposes. They think that individuals' private information should be protected [1], [2]. Additionally, safeguards for personal data protection should be in place. Privacy-preserving data mining tools have been proven and put into practice to solve this issue [3]. Security assurance solutions that have been developed based on a number of annoyances are being combined using a variety of data mining techniques. Privacy-preserving data mining (PPDM) helps to safeguard private and sensitive information for individuals. Description of architectural design of privacy preservation in data mining is shown in Figure 1. In Figure 1(a), various stages of PPDM process is shown. In order to assure information preservation, this paper aims to apply perturbation methods [4]. By substituting some alternative data that are similar to those of records with comparable non-sensitive data, all sensitive data is replaced. It can be carried out using either the distributions of sensitive data when specific non-sensitive data are present the mean of sensitive data [5], [6]. These transformations and conversions typically include numerical data. Additionally, some procedures entail simple changes [7]. By transforming user input into an improbable and unpredictable form,

---

perturbation techniques have been created as a way to ensure secrecy [8]. The data can only be modified by authorized individuals, as illustrated in the Figure 1(b). The owner then makes the data available to analysts for the data mining process [9].



Figure 1. Description of architectural design of privacy preservation in data mining (a) stages of PPDM and (b) framework of privacy preserving

Privacy and security of datasets have been the subject of extensive investigation. PPDM strategies have been proposed and used in a variety of ways. However, the majority of these strategies do not work in all situations. Mehta and Rao [10] identified existing ways from the field of natural language processing (NLP) to transform the unstructured data to a structured form. A perturbation-based method for protecting privacy in data mining was presented by Mariammal [11]. It is basis is the additive rotation strategy. In this method, the author calculated the privacy level using the variance of the initial dataset. Banu and Nagaveni [12] provided a data modification-based method to protect confidentiality in data mining process. Its foundation is the random rotation technique. They calculated the privacy level using the variance of the initial dataset. Rao *et al.* [13] demonstrated that this strategy is more efficient and accessible. They contrasted their algorithm with the perturbation strategy and demonstrated that their algorithm provided 100% data usefulness. Mary [14] asserts that the random projection strategy has a higher level of privacy than the other approaches. The photos can be very well maintained by employing RP. This method makes it possible to protect data better. It is possible to increase privacy. Ghosh *et al.* [15] gave a thorough analysis of the currently employed PPDM approaches and categorized the different data modification techniques. Javid *et al.* [16] provided a practical hybrid method for safeguarding the dataset's privacy. For numerical data, they employ geometric data perturbation, and for categorical data, they use the k anonymization technique. In their approach, they performed perturbation with randomization (intervals). Pika [17] investigated a number of data perturbation techniques in healthcare. In data perturbation methods, records' data values are changed. Their research indicates that the perturbation approach utilized to protect the confidentiality of original values.

Perturbation approach involves altering the original dataset's structure or introducing a small amount of noise to the data. By transforming user input into an improbable and unpredictable form, data perturbation may be utilized to efficiently employ PPDM. It is one of the often used techniques for protecting privacy [18]. There are several ways of perturbation. A range of methods, including as noise addition, data hiding techniques, swapping, and many more, may be used to change the information in datasets [19]. The probability distribution method and the value distortion method are two strategies of perturbation [20], [21]. In first technique, the data is immediately replaced by the distribution, however in the value distortion method, the data is directly altered either by using another randomization process or by introducing noise. Perturbation can be of three types: projection, geometric, and random perturbation [22]. In projection perturbation, modification is accomplished by changing the dimensions. Data randomly moves in this manner from high-dimensional to low-dimensional space [23]. In the geometric perturbation approach, perturbation is performed using a mixture of several techniques, including rotation transformation, translation transformation, and adding random value [24]–[26].

When publishing data, the data owner may employ a variety of privacy-preserving techniques. Owner can use various perturbation strengths to change the datasets [27]. This research offered a comparison of approaches based on random projection and principal component analysis that concurrently improve data classification accuracy while lowering the high dimension to a low dimension in order to safeguard the dataset's privacy [28], [29]. The performance of two perturbation-based privacy-preserving methods is examined in the current research. Healthcare datasets have been used to test this analysis. The following are this paper's main contributions: i) The present research provides an exhaustive analysis of PPDM techniques based on perturbation; ii) Experimental and comparative analysis of two perturbation-based privacy-preserving methods i.e., improved random projection perturbation (IRPP) and enhanced principal component analysis-based technique (EPCAT) are described in this paper; and iii) The impact of hybrid privacy preserving approaches is analyzed in this research.

The further flow of the paper is organized as: research method used in this research is presented in section 2. Section 3 provides the description of the experimental results and its comparisons with the existing work. Section 4 concludes the paper in the end.


## 2. METHOD

Numerous PPDM-related methods and techniques have previously been created and used [30]. In this article two hybrid techniques named enhanced principal component analysis-based technique and improved random projection perturbation are discussed. The efficiency of perturbation-based datasets transformation is also investigated via comparative analysis. These techniques are as follows:

### 2.1. Enhanced principal component analysis-based technique

It is a principal component analysis (PCA) and classification-based approach that protects privacy. In this method, the initial stage involves pre-processing the original data using a data filter. The filtered data is then subjected to a PCA-based modification after the data pre-processing stage. Finally, the modified data is subjected to a classification approach for data mining.

The following two phases make up the full structure of this technique:

a. Phase 1: The preservation of individual privacy in datasets is the focus of this phase. This phase consists of mainly two parts. Which are: i) The most crucial component for improving the precision and speed of the classification approach is the classification filter module (CFM). Prior to the PCA modifications of the dataset, this filter is applied to the original dataset and ii) The second module is the perturbation module, where the altered data set is once more disturbed using PCA-based transformations. Additionally examined and contrasted with the original dataset is the affected dataset's correctness.

b. Phase 2: The perturbed data set is mined after the two aforementioned modules. The "naive Bayes" approach is used as the classification method in this instance. Additionally, accuracy is calculated on the original datasets and contrasted with the accuracy of the perturbed dataset. Figure 2 shows the functional flow diagram for this model.



Figure 2. Representation of framework enhanced principal component analysis-based technique

## 2.2. Improved random projection perturbation

Random projection is a potent method that involves utilizing a random $k \times d$ matrix to project the original high d-dimensional data onto a smaller k-dimensional subspace [31]. Figure 3 provides a general overview of the design view of improved random projection perturbation method. The technique's overall structure is split into two sections.

a. Phase 1: The preservation of individual privacy in datasets is the focus of this phase. This phase consists of two parts: i) Feature selection: This module is used to choose features and improve the classification technique's accuracy. Prior to the dataset's changes using random projection, this was done to the original dataset using PCA. Prior to the random projection process and the classification phase, feature selection is used. In this paper, a feature selection method based on PCA is applied for selection of appropriate features and ii) Random projection: The perturbed data is adjusted once more in this module utilizing dimension reduction, which is accomplished with the aid of the random projection method. The datasets are distorted using the random projection technique. Additionally examined and contrasted with the original dataset is the affected dataset's correctness.

b. Phase 2: In this phase, perturbed data sets are mined using a particular classification technique. Naïve Bayes classifier is used in this instance. Additionally, several matrices are computed on the original datasets and their accuracy is contrasted with that of the perturbed dataset.



Figure 3. Architectural design of improved random projection perturbation

## 2.3. Simulation

WEKA and R-studio software are used to carry out the performance analysis of both techniques [32]. The original datasets are used in these experimental analysis together with the chosen methods to create the transformed datasets [33]. On both datasets, numerous metrics including correctly classified instances, incorrectly classified instances, true positive (TP) rate, F-measures, model building time are calculated. These metrics are used to assess how well the chosen algorithms performed on the projected dataset. The naive Bayes classification method used for implementations of the algorithms in order to evaluate the efficacy of the strategies.

### 2.3.1. Naïve Bayes

A number of classification algorithms are involved in naive Bayes classifier. These are based on Bayes' theorem. Mathematically, Bayes' theorem is stated as (1):

$$P(A\backslash B) = \frac{P(B\backslash A)\, P(A)}{P(B)} \tag{1}$$

where, $A$ and $B$ are events and $P(A\backslash B)$ is the probability of event $A$ occurring, while $B$ has already occurred, $P(B\backslash A)$ determines the probability of event $B$ occurring, while $A$ has already occurred, $P(A)$ defines the probability of event A occurring. $P(B)$ defines the Probability of event $B$ occurring.

In this performance evaluation, classification accuracy and time usage are the two key emphasis points. Comparisons are made between the performance results and the results collected from the NB. Two datasets were perturbed by the IRPP and EPCAT. The entire 10-fold cross-validation process used in the techniques. For the purpose of validation, ten samples were chosen ten times this validation process was repeated. Nine samples were collected for training in a single run. One sample was utilized to evaluate the effectiveness of the suggested method. The average accuracy of the 10 iterations was then used to represent the overall performance for a particular dataset. Figure 4 depicted the snapshots of principal components of hypothyroid dataset after performing PCA based transformations on training set and test set.

Figure 5 depicts the transformed datasets after perturbation for the cardiovascular system in Figure 5(a) and for the hypothyroid dataset in Figure 5(b). Modified datasets are shown to be more secure than the original dataset since it is challenging to access the changed data. So, confidentiality of datasets is maintained.



Figure 4. Snapshots of principal components of cardiovascular dataset on training set and test set after EPCAT perturbation technique in R-studio



(a)                                                    (b)

Figure 5. Snapshots of datasets after implementation of IRPP perturbation technique in WEKA tool (a) cardiovascular dataset and (b) hypothyroid dataset

### 2.3.2. Evaluation metrics

The effectiveness of the used techniques is assessed using various categorization measures. These include F-measures, TP rate, accuracy, runtime, and false positive (FP) rate.

- Accuracy: Accuracy can be used to evaluate a classification model. It is one aspect that may be considered to rate a classification models. Accuracy is the proportion of forecasts that our model successfully predicted. The official definition of accuracy is as (2).

$$\frac{\text{Number of correct predictions}}{\text{Total number of predictions}} \tag{2}$$

- TP rate (accuracy +): It is the fraction of accurate forecasts in positive class predictions. Its threshold is set by default to 80%. The measure is becoming better, as seen by its rising trend. Its descending trend suggests that the measure is getting worse. Data from feedback is increasingly diverging from training data. Irregular variance shows the inconsistency of feedback data. The true positive rate is calculated by (3).

$$\frac{\text{Number of True Positives}}{(\text{Number of True Positives} + \text{Number of False Negatives})} \tag{3}$$

- FP rate (accuracy-): Accuracy of machine learning algorithms may be evaluated using a statistic called the False Positive Rate. The false positive rate is determined as (4).

$$\frac{FP}{(FP+TP)} \tag{4}$$

- F-measures: It is combined measure for precision and recall metrics. It provides a single score that balances both the concerns of precision and recall in one number. It is calculated as (5).

$$\frac{2 * \text{Precision} * \text{Recall}}{(\text{Precision} + \text{Recall})} \tag{5}$$

## 3. RESULTS AND DISCUSSION

This research analyses the impact of hybrid perturbation based PPDM techniques in data mining process. For this purpose, IRPP and EPCAT are selected. Several tests were run on data sets of two different sizes, and the associated outcomes were seen. The results of the experiments demonstrate that the both hybrid techniques IRPP and EPCAT perform better due to their greater accuracy, TP rate, FP rate, F-measures, and run duration values.

- Datasets: This research is implemented and experimented on two datasets i.e., cardiovascular dataset and the hypothyroid dataset. The cardiovascular dataset consists of 70k instances and 13 attributes. The hypothyroid dataset consists of 7,200 instances and 21 attributes [34].

The effectiveness of both method on cardiovascular datasets is shown in Table 1. On the provided training datasets, the metrics accuracy, TP rate, FP rate, F-measures, and run time are computed. It is clearly shown in the Table 1 and can be noticed that the strategies produced superior results than the conventional model of categorization in all regards.

Table 1. Performance measure of naïve Bayes classifier on cardiovascular dataset

| Accuracy measurement | Original dataset | IRPP based perturbed dataset | EPCAT based perturbed dataset |
|---|---|---|---|
| Time taken to build (sec) | 0.53 | 0.34 | 0.25 |
| Correctly classified instances (%) | 58.85 | 59.87 | 59.65 |
| Incorrectly classified instances (%) | 41.15 | 40.13 | 40.31 |
| TP rate (%) | 58.90 | 59.70 | 59.16 |
| F-measures (%) | 54.40 | 56.50 | 56.14 |

The performance of IRPP and EPCAT privacy-preserving algorithms to the conventional classification algorithms on cardiovascular datasets is shown in Figure 6. It displays the efficiency of the random projection-based privacy-preserving and PCA-based privacy-preserving method to the traditional classification model on cardiovascular datasets. As shown in the figure, it is observed that the IRPP method has better accuracy measures than the conventional classification algorithms and PCA-based privacy-preserving method. The effectiveness of both methods to the conventional classification model on hypothyroid dataset is shown in Table 2.

On provided training datasets, the metrics accuracy, TP rate, FP rate, F-measures, and run time are calculated. It is well depicted in the in the table, and it is easy to see that the IRPP approach yields better results overall than the conventional model of classification. On hypothyroid datasets, Figure 7 compares the efficiency of the privacy-preserving algorithms IRPP and EPCAT to the traditional classification techniques. It demonstrates the effectiveness of the privacy-preserving random projection and PCA methods in comparison to the conventional classification model on hypothyroid datasets. As depicted in the Figure 7, it can be seen that the IRPP approach outperforms both the PCA-based privacy-preserving method and traditional classification algorithms in terms of accuracy measurements.

Figure 6. Performance analysis IRPP and EPCAT to the conventional model on cardiovascular dataset

Table 2. Performance measure of naïve Bayes classifier on hypothyroid dataset

| Accuracy measurement | Original dataset | IRPP based perturbed dataset | EPCAT based perturbed dataset |
|---|---|---|---|
| Correctly classified instances (%) | 71.67 | 75.67 | 74.08 |
| Incorrectly classified instances (%) | 28.32 | 24.32 | 25.91 |
| TP rate (%) | 71.70 | 75.13 | 74.10 |
| F-Measures (%) | 70.80 | 73.54 | 73.48 |



Figure 7. Performance Analysis IRPP and EPCAT to the conventional model on hypothyroid dataset using naïve Bayes classifier

## 4. CONCLUSION

Developing algorithms that can conceal or provide privacy to some sensitive information is the fundamental goal of privacy preservation in data mining operations. PPDM techniques are essential in order to stop profiteers from gaining unwanted access. However, data mining accuracy and privacy conflict. In this context, this paper has analyzed the impact of PPDM techniques based on perturbation to datasets. This article provides a brief overview of some privacy techniques, namely PCA-based perturbation, and random projection-based perturbation, and analyzes their competencies and differences in different scenarios. The effectiveness of these hybrid techniques has been tested in classification algorithms naive Bayes classifiers. For implementation purpose two datasets cardiovascular and hypothyroid datasets have been selected. It has been found that IRPP privacy-preserving approach and enhanced principal component analysis-based technique EPCAT are more effective than traditional technique. The perturbed datasets are more privacy preserved. In cardiovascular dataset's case, the perturbed datasets outperform the original dataset in terms of runtime, accuracy, TP rate, and F-measurer. In hypothyroid dataset's case, implementation results on all measurements are better or almost identical to the previous approach model. Therefore, it is noticed that the datasets that are altered using hybrid privacy preserving approaches are more secure and efficient than the original datasets.

## REFERENCES

[1] A. Altalhi, M. AL-Saedi, H. Alsuwat, and E. Alsuwat, "Privacy-preserving in the context of data mining and deep learning," *International Journal of Computer Science and Network Security*, vol. 21, no. 6, pp. 137–142, 2021.

[2] M. T. Andavan and N. Vairaperumal, "Privacy protection domain-user integra tag deduplication in cloud data server," *International Journal of Electrical and Computer Engineering (IJECE)*, vol. 12, no. 4, pp. 4155–4163, Aug. 2022, doi: 10.11591/ijece.v12i4.pp4155-4163.

[3] P. Gulia and Hemlata, "Privacy preserving data mining of vertically partitioned data in distributed environment-an experimental analysis," *Journal of Theoretical and Applied Information Technology*, vol. 96, no. 10, pp. 2973–2987, 2018.

[4] R. Ratra and P. Gulia, "Privacy preserving data mining: techniques and algorithms," *International Journal of Engineering Trends and Technology*, vol. 68, no. 11, pp. 56–62, Nov. 2020, doi: 10.14445/22315381/IJETT-V68I11P207.

[5] H. Kargupta, S. Datta, Q. Wang, and K. Sivakumar, "Random-data perturbation techniques and privacy-preserving data mining," *Knowledge and Information Systems*, vol. 7, no. 4, pp. 387–414, May 2005, doi: 10.1007/s10115-004-0173-6.

[6] Kun Liu, H. Kargupta, and J. Ryan, "Random projection-based multiplicative data perturbation for privacy preserving distributed data mining," *IEEE Transactions on Knowledge and Data Engineering*, vol. 18, no. 1, pp. 92–106, Jan. 2006, doi: 10.1109/TKDE.2006.14.

[7] G. N. Devi and K. Manikandan, "Improved perturbation technique privacy-preserving rotation-based condensation algorithm for privacy preserving in big data stream using internet of things," *Transactions on Emerging Telecommunications Technologies*, vol. 31, no. 12, Dec. 2020, doi: 10.1002/ett.3970.

[8] K. M. Chong, "Privacy-preserving healthcare informatics: a review," *ITM Web of Conferences*, vol. 36, Jan. 2021, doi: 10.1051/itmconf/20213604005.

[9] M. Dabhade and J. J. Hilda, "Privacy preserving in data mining using data perturbation and classification method," *Emerging trends in Computer Engineering and Research*, vol. 8, no. 2, pp. 346–352, 2017.

[10] B. B. Mehta and U. P. Rao, "Improved l-diversity: scalable anonymization approach for privacy preserving big data publishing," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 4, pp. 1423–1430, Apr. 2022, doi: 10.1016/j.jksuci.2019.08.006.

[11] S. Mariammal, "An additive rotational perturbation technique for privacy preserving data mining," *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, vol. 12, no. 9, pp. 2675–2681, 2021.

[12] R. V. Banu and N. Nagaveni, "Preservation of data privacy using PCA based transformation," in *2009 International Conference on Advances in Recent Technologies in Communication and Computing*, 2009, pp. 439–443, doi: 10.1109/ARTCom.2009.159.

[13] P. R. M. Rao, S. M. Krishna, and A. P. S. Kumar, "Novel algorithm for efficient privacy preservation in data analytics," *Indian Journal of Science and Technology*, vol. 14, no. 6, pp. 519–526, Feb. 2021, doi: 10.17485/IJST/v14i6.1773.

[14] Av. Mary, "A random projection approach to secure medical images," *International Journal of Advanced Research*, vol. 7, no. 3, pp. 1298–1301, Mar. 2019, doi: 10.21474/IJAR01/8763.

[15] S. Ghosh, S. Sadhu, S. Biswas, D. Sarkar, and P. P. Sarkar, "A comparison between different classifiers for tennis match result prediction," *Malaysian Journal of Computer Science*, vol. 32, no. 2, pp. 97–111, Apr. 2019, doi: 10.22452/mjcs.vol32no2.2.

[16] T. Javid, M. K. Gupta, and A. Gupta, "A hybrid-security model for privacy-enhanced distributed data mining," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 6, pp. 3602–3614, Jun. 2022, doi: 10.1016/j.jksuci.2020.06.010.

[17] A. Pika, M. T. Wynn, S. Budiono, A. H. M. Ter Hofstede, W. M. P. van der Aalst, and H. A. Reijers, "Privacy-preserving process mining in healthcare," *International Journal of Environmental Research and Public Health*, vol. 17, no. 5, p. 1612, Mar. 2020, doi: 10.3390/ijerph17051612.

[18] R. Ratra, P. Gulia, and N. S. Gill, "Evaluation of re-identification risk using anonymization and differential privacy in healthcare," *International Journal of Advanced Computer Science and Applications*, vol. 13, no. 2, pp. 563–570, 2022, doi: 10.14569/IJACSA.2022.0130266.

[19] T. A. Adesuyi and B. M. Kim, "A layer-wise Perturbation based Privacy Preserving Deep Neural Networks," in *2019 International Conference on Artificial Intelligence in Information and Communication (ICAIIC)*, Feb. 2019, pp. 389–394, doi: 10.1109/ICAIIC.2019.8669014.

[20] R. Ratra, P. Gulia, N. S. Gill, and J. M. Chatterjee, "Big data privacy preservation using principal component analysis and random projection in healthcare," *Mathematical Problems in Engineering*, vol. 2022, pp. 1–12, Aug. 2022, doi: 10.1155/2022/6402274.

[21] A. Amkor, K. Maaider, and N. El Barbri, "An evaluation of machine learning algorithms coupled to an electronic olfactory system: a study of the mint case," *International Journal of Electrical and Computer Engineering (IJECE)*, vol. 12, no. 4, pp. 4335–4344, Aug. 2022, doi: 10.11591/ijece.v12i4.pp4335-4344.

[22] E. O. Abiodun, A. Alabdulatif, O. I. Abiodun, M. Alawida, A. Alabdulatif, and R. S. Alkhawaldeh, "A systematic review of emerging feature selection optimization methods for optimal text classification: the present state and prospective opportunities," *Neural Computing and Applications*, vol. 33, no. 22, pp. 15091–15118, Nov. 2021, doi: 10.1007/s00521-021-06406-8.

[23] S. A. Abdelhameed, S. M. Moussa, N. L. Badr, and M. E. Khalifa, "The generic framework of privacy preserving data mining phases: challenges & future directions," in *2021 Tenth International Conference on Intelligent Computing and Information Systems (ICICIS)*, Dec. 2021, pp. 341–347, doi: 10.1109/ICICIS52592.2021.9694174.

[24] X. Fan, G. Wang, K. Chen, X. He, and W. Xu, "PPCA: privacy-preserving principal component analysis using secure multiparty computation (MPC)," *arXiv:2105.07612*, May 2021.

[25] Z. Chen and K. Omote, "A privacy preserving scheme with dimensionality reduction for distributed machine learning," in *2021 16th Asia Joint Conference on Information Security (AsiaJCIS)*, Aug. 2021, pp. 45–50, doi: 10.1109/AsiaJCIS53848.2021.00017.

[26] M. A. Abdo, A. A. Abdel-Hamid, and H. A. Elzouka, "A cloud-based mobile healthcare monitoring framework with location privacy preservation," in *2020 International Conference on Innovation and Intelligence for Informatics, Computing and Technologies (3ICT)*, Dec. 2020, pp. 1–8, doi: 10.1109/3ICT51146.2020.9311999.

[27] M. B. Malik, M. A. Ghazi, and R. Ali, "Privacy preserving data mining techniques: current scenario and future prospects," in *2012 Third International Conference on Computer and Communication Technology*, Nov. 2012, pp. 26–32, doi: 10.1109/ICCCT.2012.15.

[28] P. H. Li, T. Lee, and H. Y. Youn, "Dimensionality reduction with sparse locality for principal component analysis," *Mathematical Problems in Engineering*, vol. 2020, pp. 1–12, May 2020, doi: 10.1155/2020/9723279.

[29] T. Meurers, R. Bild, K.-M. Do, and F. Prasser, "A scalable software solution for anonymizing high-dimensional biomedical data," *GigaScience*, vol. 10, no. 10, Oct. 2021, doi: 10.1093/gigascience/giab068.

[30] P. Churi, A. Pawar, and A.-J. Moreno-Guerrero, "A comprehensive survey on data utility and privacy: taking Indian healthcare system as a potential case study," *Inventions*, vol. 6, no. 3, Jun. 2021, doi: 10.3390/inventions6030045.

[31] T. Bianchi, V. Bioglio, and E. Magli, "Analysis of one-time random projections for privacy preserving compressed sensing," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 2, pp. 313–327, 2016, doi: 10.1109/TIFS.2015.2493982.

[32] R. Ratra and P. Gulia, "Experimental evaluation of open source data mining tools (WEKA and orange)," *International Journal of Engineering Trends and Technology*, vol. 68, no. 8, pp. 30–35, Aug. 2020, doi: 10.14445/22315381/IJETT-V68I8P206S.

[33] S. K. David, M. Rafiullah, and K. Siddiqui, "Comparison of different machine learning techniques to predict diabetic kidney disease," *Journal of Healthcare Engineering*, vol. 2022, pp. 1–9, Apr. 2022, doi: 10.1155/2022/7378307.

[34] "Find open datasets and machine learning projects," *Kaggle*. Accessed: May 31, 2022. [Online]. Available: https://www.kaggle.com/datasets.

## BIOGRAPHIES OF AUTHORS

**Ritu Ratra** 🔾 🈳 SC ▷ received the MCA degrees from Maharshi Dayanand University, Rohtak, and Haryana, India. Currently, she is pursuing her Ph.D. in Computer Science at the Department of Computer Science and Applications, Maharshi Dayanand University. Rohtak, Haryana, India. Previously, she had worked as an Assistance Professor at Sh. Lal Nath Hindu College, Rohtak, Haryana, India affiliated with MD University, Rohtak, and Haryana, India for 12.5 years. Her research areas include machine learning, data mining, big data, and artificial intelligent. She has authored or coauthored many refereed journal and conference papers. She can be contacted at email: ritu.rs.dcsa@mdurohtak.ac.in.

**Preeti Gulia** 🔾 🈳 SC ▷ received Ph.D. degree in computer science in 2013. She is currently working as Associate Professor at the Department of Computer Science and Applications, M.D. University, Rohtak, Haryana, India. She is serving the Department since 2009. She has published more than 65 research papers and articles in journals and conferences of National/International repute including ACM, and Scopus. Her area of research includes data mining, big data, machine learning, deep learning, IoT, and software engineering. She is an active professional member of IAENG, CSI, and ACM. She is also serving as editorial board member active reviewer of international/national journals. She has guided four research scholars as well as guiding six Ph.D. research scholars from various research areas. She can be contacted at email: preeti@mdurohtak.ac.in.

**Nasib Singh Gill** 🔾 🈳 SC ▷ holds Post-Doctoral research in Computer Science at Brunel University, West London during 2001-2002 and Ph.D. in Computer Science in 1996. He is a recipient of the Commonwealth Fellowship Award of the British Government for the Year 2001. Besides, he also has earned his MBA degree. He is currently Head, Department of Computer Science and Applications, M. D. University, Rohtak, India. He is also working as Director, Directorate of Distance Education as well as Director of Digital Learning Centre, M. D. University, Rohtak, Haryana. He is an active professional member of IETE, IAENG, and CSI. He has published more than 304 research papers and authored 5 popular books He has guided so far 12 Ph.D. scholars as well as guiding about 5 more scholars. His research interests primarily include-IoT, machine and deep learning, information and network security, data mining and data warehousing, NLP, and measurement of component-based systems. He can be contacted at email: nasib.gill@mdurohtak.ac.in.