

A Privacy-Preserving Data Mining Approach in Multi-Dimensional Data Set based on the Random and Cumulative Integrated Noise

Shailesh Kumar Vyas*¹, Swapnil Karmore², Pinkal Jain³

Submitted: 10/12/2023 Revised: 14/01/2024 Accepted: 31/01/2024

Abstract: Since the emergence of Big Data, data mining has changed for the better. We now have unprecedented opportunities to discover new knowledge and support decision-making. While Big Data methods are gaining traction, privacy is becoming a major concern. In order to overcome these difficulties, we employ new data perturbation approaches based on random translation and random projections, as well as additive noise. In addition, we assess the performance of different data perturbation modes and their corresponding attack modes using input-output maximum a posteriori (MAP) attacks. Firstly, we evaluate random projections for online classifications. Two data perturbation modalities are assessed: random translation (RT) and random projection (R). Independent noise (RPIN), or cumulative noise (RCPN) are also assessed. Our results show that a combination of 2 MAP attacks (MAX (A-RT, A-RCPN-1)) vs. RCPN method is the most efficient. As the data record moves away from the known data record, the attack becomes less efficient, indicating that over time, RCPN gradually improves data privacy. Thanks to our work, we can apply these perturbation techniques to more than just classification tasks. We can apply them to cluster, anomaly detection and regression, which opens up new research directions. We are also exploring privacy preservation techniques that are tailored to the streaming nature of real-world data sets to improve privacy for nominal data.

Keywords: Cumulative noise, data mining, data perturbation, privacy preservation, RCPN

1. Introduction

Big Data has popularized the field of data mining, which has led to the discovery of new knowledge and the provision of decision support. However, there are many challenges in the implementation of big data methods that need to be addressed before these methods can be applied to real-world situations. An important concern for users is privacy protection in data mining. A majority of users are hesitant to share private information which restricts the useful data to be available for mining purposes. The Netflix Prize dataset, which proved that users may be recognized using their publicly available movie ratings submitted on the Internet Movie Database, is an example of placing insufficient importance to privacy on data exchange platforms [26].

To avoid such issues in the future, the community of data mining developed several methods for PPDM [22] using which privacy of users was protected while using their data for analysis purposes. PPDM relates to privacy preservation during data mining, whereas PPDP generates a clean dataset which can be used and shared without exposing sensitive information of users. These methods convert sensitive information to a format that protects users' privacy and that can be used to build models having an accuracy which is

comparable to that of models built using the original information.

There are many limitations of data streams, including humongous information arising from continuously generated and transmitted data, and pressing require to quickly build models to analyze this data by recognizing underlying patterns. In such a fast-paced data generation environment, any privacy-preservation method will need to make sure that the method is capable of handling the speed at which data is generated and can perform efficiently. When data streams have high data generation rates, the overheads responsible for privacy preservation should be minimal so that transformations are implemented before the next data arrives; otherwise, it might lead to load shedding and loss of information. If this data is useful only for specific time following its arrival, then the method should be able to exert a high level of computation to ensure that sensitive record is not exposed in that time frame.

Several PPDM and PPDP methods have high overheads for implementation of privacy preservation protocols and are therefore not suitable for data mining [17]. As a result, this has led to the development of specific privacy preservation strategies for data stream mining. An important consideration for privacy preservation methods is that there needs to be a balance between data privacy and utility. The greater the data privacy, the lesser it's utility for analysis thereby lowering the accuracy of the model. Despite this, it is extremely important to ensure privacy when highly sensitive data is being handled. The accuracy of the model is useful only if it ensures privacy of data. Therefore, current

¹ G H Raisoni University, Saikheda, MP-480106
<https://orcid.org/0009-0003-1262-7353>

² G H Raisoni University, Saikheda, MP-480106
<https://orcid.org/0000-0003-2068-4043>

³ Amity University, Gwalior, MP-474020
<https://orcid.org/0000-0001-8002-320X>

* shailesh.pk.29@email.com

research in this field is concerned with optimizing the balance between data privacy and utility. Also, newly developed privacy preservation methods are constantly under attack [27]. In our study, we have explored a combination of various mechanisms of data perturbation for either PPDM or PPDP for classifying data streams online with the objective of striking a balance between model accuracy and data privacy. The main contributions of our study are:

We applied relatively novel data perturbation methods that utilize the concepts of randomized translation, randomized projection, and two kinds of added noise. We modified the input-output Maxim on randomized basis for addition of mixed noise and randomized translation. We have evaluated the various data perturbation methods along with their corresponding attack methods. First to examine the performance of randomized projection for the aim of web-based author classification.

2. Literature Review

Privacy is defined as concealing the data pertaining to a person's usage of a website or an application as it may have negative consequences for that person [12]. As a result, privacy preservation refers to the implementation of strategies that protects the browsing data of people from access by unauthorized persons. Several studies have focused on this aspect and many efforts have been made in the field of PPDM. Perturbation algorithm was proposed for a distributed scenario known as DISTPAB for privacy preservation of data partitioned horizontally [2, 3]. In this method, the asymmetrical patterns of resources are used for distributing the privacy preservation tasks and removing the computational hurdles. It is suitable for both resource-constraint devices as well as high-performing systems. An approach for privacy preservation based on distance matrix was developed to protect shared data among different organizations [7]. The original data's distance matrix was calculated using the shared data and this distance matrix was shared among organizations to mimic the original data.

The Reversible Privacy Contrast Mapping (RPCM) algorithm was suggested that distorts and restores data using reversible strategies for concealing data [27]. Experiments have indicated that RPCM does not lead to data loss when distorting data. Shan et al. (2020) used range noise to propose a novel random perturbation strategy known as the Range Noise Perturbation (RNP) method. This method performs data perturbation by: (1) selecting a set of data characteristics, (2) selecting a value for range noise, (3) selecting the noise after calculating the range, and (4) updating the dataset with the help of the perturbed data characteristics.

A data perturbation strategy based on a Bayesian model was discussed [21]. This method does not use an algorithm and

is best for classification. This shows that the classification algorithms that are currently in use can directly use the perturbed data ensuring privacy preservation. A clustering algorithm based on semi-supervised schema was proposed for privacy preservation that makes use of small supervised information [16]. This algorithm first computes the metric for cluster recognition with the help of convex optimization. Following this, the learned metric is used to impose a multiplicative alternation on actual data set which is used to alter the dimensionality of actual data set. This not only aids privacy preservation but also ensures that the original features of the data are conserved.

A framework for spatial data transformation was proposed which was known as rotation-based transformation [28]. In this method, operations of mathematical foundation, clustering, and other CPU-intensive operations were performed efficiently. A method for data perturbation using min-max normalization has been suggested [20]. It involves the selection of sensitive characteristics from the dataset after which the min-max normalization is applied on the selected characteristics for perturbation. These perturbed features are then used for analysis purposes following integration with the original dataset. A method based on 3D rotation transformation was discussed that performs the rotation of data characteristics in groups of three each [8]. The angle of rotation is chosen based on the variance of the characteristics that threshold to a bias factor was described [21]. A technique has been developed that is based on rotation and condensation and it is also known as the $P^2 RoC AI$. This algorithm uses the condensation and rotation characteristics to preserve the privacy and it maintains the actual features of the data in data streams. PABIDOT, a non-reversible method was developed for data perturbation, which is based on optimal geometric transformations [2, 3]. It is useful for privacy preservation of big data. It functions by using multi-dimensional geometrical alternations, translations, reflections, and rotations for perturbing data following which, the data is subjected to random data set generation along with random tuple shift out. Two methods have been developed to compute the data perturbation that leads to the privacy preservation in mining [10]. These two methods are based on random translation, random projection, and there is a couple of techniques for noise integration, one of which is independent and generated for individual purposes and the other is accumulative and generated across the data's lifespan.

SEAL, an efficient and secure method was proposed for data perturbation, which is based on local differential privacy [6]. This method was developed using Chebyshev interpolation and Laplacian noise. There is an ideal for method for privacy preservation in data mining which is generated through smart cyber-physical systems. A 2-stage schema was proposed that is based on RG (Repeated Gompertz) + RP (Random Projection) [25]. RG is a non-

linear function that is used in the first step for making changes to the data. Then, the RP matrix is used to reduce the data dimensions while preserving distances. RG is designed such that it prevents MAP estimation attacks and RP prevents ICA attacks thereby confirming clustering performance.

Data perturbation using a fuzzy data transformation approach is performed [7]. In this, the fuzzifier is used to convert the crisp values into linguistic values and an inference engine is used to compare these values with the fuzzy rules previously defined. Then, a defuzzifier is used to convert the linguistic values back into the crisp values which are used as the data which has been perturbed.

3. Models of Data Perturbation

This section describes the models of data perturbation, specifically the foundational model that is based on translation and RP, the model with cumulative additive noise, and the model with independent additive noise. For all models, the features of the data and columns represent the records. When a data perturbation method is applied, X is converted into a perturbed dataset which is represented as $k \times n$ matrix Y . It should be noted that the number of columns remains the same between the two matrices while the number of rows may decrease ($k \leq m$).

3.1. Random Projection Techniques

This model can be represented as:

$$Y = \frac{1}{\sqrt{k}\sigma_r}RX \quad (1)$$

Where, Y is the matrix multiplication,

$R = k \times m$ (matrix).

When the projection is multiplied by $\frac{1}{\sqrt{k}\sigma_r}$, the values in the columns are preserved when horizontally-arranged values in datasets are perturbed using the same R . These datasets typically represent different records of data having the same features. This is not mandatory when only one dataset is available because when the scale of the dataset is changed, it will not have any effect on data mining.

An approach was discussed on the basis of RP [18], according to which it is possible to reduce a dataset having s records to $O\left(\frac{\log s}{\epsilon^2}\right)$ dimensions while maintaining the pairwise distances intact (ϵ represents error). As a result, RP perturbs data while maintaining distances. With a reduction in k , there is a corresponding exponential increase in the error of pairwise distances [5, 6]; however, the resistance of the system to attacks also increases [16].

3.1.1. Using random translation for resisting rotation centre attacks

The reduction in dimensionality brought about by RP may reduce the vulnerability of the data to attacks; however,

further reduction of the vulnerability may be brought about by applying random translation during the modification process. Therefore, the modified model can be extended to the following equation:

$$Y = \frac{1}{\sqrt{k}\sigma_r}RX + \varphi \quad (2)$$

Where every column of φ is the same. Furthermore, every element in each of the rows may be a positive or a negative number that is extracted from value (feature range F_i) and max value of translation (equal to twice the value of the range).

$$\varphi_{*,i} = \varphi_{*,j}, 1 \leq i < j \leq n \quad (3)$$

$$\varphi = \begin{pmatrix} \varphi_{i,j} \\ \vdots \\ \varphi_{k,n} \end{pmatrix} = \beta(-1,1) \times \mu(R(F_i), 2R(F_i)) \quad (4)$$

Where, $R(F_i)$ represents the properties range

In general, most data mining tasks are not affected by applying constant random translation to the data records. Despite this, it is still useful because attackers can only account for this translation by sacrificing a pair of input-output records.

Therefore, this foundational model of data perturbation involves both RP as well as random translation.

3.2. Random projection combined with independent noise

The RP method can be further strengthened by using additive noise in two different forms to balance between privacy preservation and accuracy. One of the ways in which this can be done to add a noise to each entry in the records that has undergone perturbation. In this way, the foundational model can be extended by including independent noise which is known as RPIN. It is important that the variance of this independent noise is proportional to each feature's (F_i) range and is represented as follows:

$$Y = \frac{1}{\sqrt{k}\sigma_r}RX + \varphi + \Delta \quad (5)$$

$$\Delta = \begin{pmatrix} \delta_{i,j} \\ \vdots \\ \delta_{k,n} \end{pmatrix} = N(0, \sigma_\delta^2 \cdot R(F_i)) \quad (6)$$

3.3. Random projection combined with cumulative noise

This approach is especially suited for data mining and represents a modification of RPIN, where cumulative noise is used instead of independent noise and the model is called RPCN. Similar to RPIN, each data record is integrated with Gaussian values and every subsequent data record is integrated with a random value ($\gamma_{i,j}$) in such a way that the noise accumulates in the dataset. This can be denoted by:

$$Y = \frac{1}{\sqrt{k}\sigma_r}RX + \varphi + \Gamma \quad (7)$$

$$\Gamma = \begin{pmatrix} \gamma_{i,j} \\ \vdots \\ \gamma_{k,n} \end{pmatrix} = \begin{cases} N(0, \sigma_\gamma^2 \cdot R(F_i)) & j = 1 \\ N(0, \sigma_\gamma^2 \cdot R(F_i)) + \gamma_{i,j-1} & j > 1 \end{cases} \quad (8)$$

Where, i represents the data feature, and j represents the index of the data record.

As shown in Figure 1, the successive values that are present in each row of the noise (Γ) is in the form of a Gaussian random walk which is important to protect the data against input-output attacks. When attackers try to access data that is far away from the known data, they will encounter a progressive increase in noise levels. Also, cumulative noise with a small value of σ_γ is almost the same as independent noise with a large value of σ_δ . So, when a small value of σ_γ is used, it has a minimal effect on the pairwise distances between the records of the data stream which increases the accuracy of data mining. The Gaussian random walk starts to exhibit a sustained gradual motion over time which can represent concept drift, a phenomenon that a lot of data mining algorithms tend to adapt to with time [1]. As a result, the privacy benefit of RPCN is similar to that of RPIN; however, it has a reduced impact on the accuracy of the algorithms.

3.4. Comparison between independent noise and cumulative noise

To compare independent and cumulative noise, it is important to first establish the relation among parameters of σ_δ in the same levels of noise. We can represent this relationship in the form of simple equations by assuming $R(F_i) = 1$, which indicates that the features of the dataset are min-max normalized. Thus, the Gaussian distribution defined by γ_{ij} can be simplified to $N(0, \sigma_\gamma^2)$ and the Gaussian distribution defined by δ_{ij} can be simplified to $N(0, \sigma_\delta^2)$.

The following half-normal distribution:

$$|N(0, \sigma_\delta^2)|$$

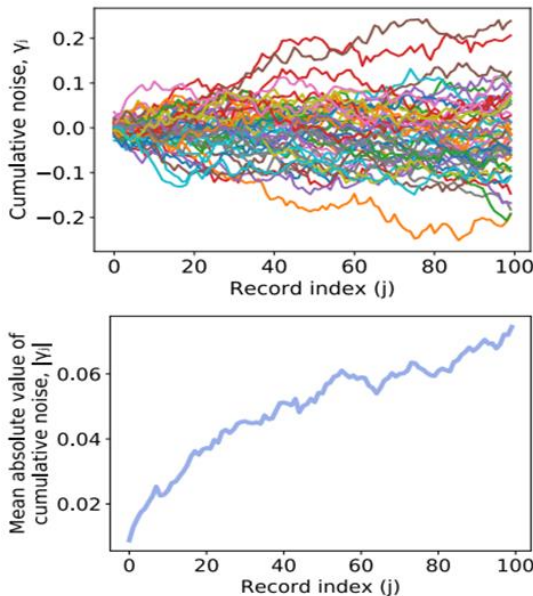


Fig 1: Random walks of cumulative noise (50 simulations); top - $\sigma_\gamma = 0.01$, bottom – average values of the absolute noise across simulations

The average of the above equation can be expressed as:

$$\sigma_\delta \cdot \sqrt{\frac{2}{\pi}}$$

Hence, an estimate of the independent noise value (E_δ) which is integrated with n number of records can be expressed as:

$$E_\delta = \sum_{i=1}^n (\sigma_\delta \cdot \sqrt{\frac{2}{\pi}}) \quad (9)$$

$$E_\delta = n \cdot \sigma_\delta \cdot \sqrt{\frac{2}{\pi}} \quad (10)$$

Similarly, we can calculate the cumulative noise (E_γ) which is integrated with n number of records; however, we must consider that the variance of the cumulative noise that is integrated with each data record increases proportionally for all subsequent data records as follows:

$$E_\gamma = \sum_{i=1}^n (\sqrt{i} \cdot \sigma_\gamma \cdot \sqrt{\frac{2}{\pi}}) \quad (11)$$

$$E_\gamma = \sum_{i=1}^n (\sqrt{i} \cdot \sigma_\gamma \cdot \sqrt{\frac{2}{\pi}}) \quad (12)$$

$$E_\gamma = \sigma_\gamma \cdot \sqrt{\frac{2}{\pi}} \cdot \sum_{i=1}^n \sqrt{i} \quad (13)$$

Thus, from above equations

$$E_\delta = E_\gamma \quad (14)$$

$$n \cdot \sigma_\delta \cdot \sqrt{\frac{2}{\pi}} = \sigma_\gamma \cdot \sqrt{\frac{2}{\pi}} \cdot \sum_{i=1}^n \sqrt{i} \quad (15)$$

$$n \cdot \sigma_\delta = \sigma_\gamma \cdot \sum_{i=1}^n \sqrt{i} \quad (16)$$

$$\sigma_\gamma = \sigma_\delta \cdot \frac{n}{\sum_{i=1}^n \sqrt{i}} \quad (17)$$

This relationship has been shown in Figure 2 where the expected noise value added to each data record $\sigma_\delta = 0.11$ becomes equivalent to noise ($\sigma_\delta \sim 0.0047$) is represented as a graph.

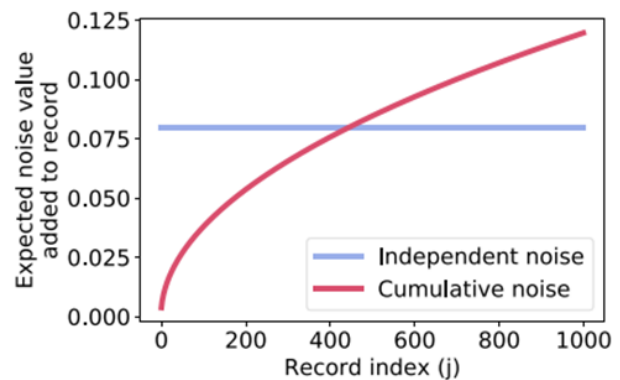


Fig 2: A graph comparing different noise when both are equal

3.5. Interpolation of cumulative noise between known data points

Consider that cumulative noise is integrated with two data records – γ_a and γ_c . Then, we can represent the cumulative noise that is integrated with a data record between these two records as γ_b , where $a < b < c$. We can estimate this cumulative noise by multiplying as follows:

$$P(\gamma_b|\gamma_a) = N(\gamma_a, (b-a)\sigma_\gamma^2) \quad (18)$$

$$P(\gamma_b|\gamma_c) = N(\gamma_c, (c-b)\sigma_\gamma^2) \quad (19)$$

$$P(\gamma_b|\gamma_a\gamma_c) = N(\mu, \sigma^2) \quad (20)$$

$$\mu = \frac{\gamma_a(c-b)\sigma_\gamma^2 + \gamma_c(b-a)\sigma_\gamma^2}{(c-b)\sigma_\gamma^2 + (b-a)\sigma_\gamma^2} \quad (21)$$

$$\mu = \frac{\gamma_a(c-b) + \gamma_c(b-a)}{c-a} \quad (22)$$

$$\sigma^2 = \frac{(b-a)\sigma_\gamma^2(c-b)\sigma_\gamma^2}{(b-a)\sigma_\gamma^2 + (c-b)\sigma_\gamma^2} \quad (23)$$

$$\sigma^2 = \frac{(b-a)(c-b)\sigma_\gamma^2}{c-a} \quad (24)$$

It must be noted that as the value of b becomes closer to the value of c or a , there is a decrease in the effective variance, because when $b = \frac{c-a}{2}$, then the value of $(b-a)(c-b)$ becomes maximum.

They may not be able to determine the value of the cumulative noise. Consider that the perturbation applied on the data perfectly preserves the pairwise distance between the data records such that $|x_c - x_a| = |y_c - y_a|$ (because pairwise distance is preserved by RP), then represented as:

$$\|\gamma_c - \gamma_a\| = |(\|y_c - y_a\| - \|x_c - x_a\|)| \quad (25)$$

However, the difference vector $\gamma_c - \gamma_a$ cannot be computed.

3.6. Efficiency of data perturbation

The performance is not significantly affected by extra operations that are implemented for additive noise and random translation. For RPIN, the noise that is generated and the random translation need to be added to each data feature (k) that is present in the projected data record in a way that the total computational complexity is proportional to $O(km)$. In case of RPCN, the noise that is generated can be integrated with the random translation following which it is integrated with the records providing the same computational complexity as provided by the noise.

4. Known Input-Output Attacks

There are several input-output attacks that can cross the barrier of the perturbation methods that have been described above. However, the prerequisite for these attacks is that the attacker should know a subset of the input data records along with their corresponding perturbed data records. This is possible when every data record represents an individual so

that a privacy breach of other data records presents in the stream.

Let us consider that the input data stream is X , and the known data (p) present in columns is represented as X_p . Similarly, consider that the output data stream is Y , and the corresponding columns of this stream are represented as Y_p . The known data will be used by the attacker to breach the privacy of other data ($y_i \in \frac{Y}{Y_p}$) so that data record (x_i) is generated. According to existing conventions proposed by [19], it is likely that the data perturbation method applied to the input record is known to the attacker along with the variances (σ_r^2 , σ_δ^2 , and σ_γ^2). The attacker may either reveal this information or it may be shared by external organizations that use the same approach for perturbation to share and merge data. As stated previously, we consider that the features have undergone min-max normalization so that the expressions of privacy attacks can be simplified.

4.1. Notations

The following notations will be followed in the rest of the paper:

- $[A, a]$ stands for the matrix that is generated by integrating a vector (a) with the matrix (A) as a new column. We consider that the columns in the matrix are sorted based on their order of addition especially when they are related to the data records.
- \bar{A} stands for the matrix that is generated by combining all columns of the matrix $m \times n$ such that they form a single column vector having length mn .
- δ_i is the independent noise vector and γ_i is the cumulative noise vector, both of which form a part of the modified data record (y_i).

4.2 Input-output MAP attack on RP

This MAP attack can perform a privacy breach of datasets having known pairs of inputs and outputs that have undergone RP. For this attack, X_p needs to have a complete column rank, meaning that all the columns need to be linearly independent of each other. Furthermore, the target x_i also needs to be linearly independent of all the columns of X_p . Therefore, the value of p must be lower than that of m .

In order to carry out an attack, \hat{x}_i is estimated as \hat{x} in order to maximize the probability of an RP of $[X_p, \hat{x}]$ resulting in $[Y_p, y_i]$.

$$\hat{x}_i = \arg \sup_{\hat{x}} \phi_r(\overline{[Y_p, y_i]}) \quad (26)$$

Where, ϕ_r and $\hat{x} \in R^m$ l.i. X_p denote the probability density function of the RP (Y_p) which is represented as follows:

$$\overline{[Y_p, y_i]} = \frac{i}{\sqrt{k}\sigma_r} \overline{R[X_p, \hat{x}]} \quad (27)$$

Φ_r is distributed as a Gaussian distribution comprising of dimensions having a block-diagonal covariance matrix and a zero mean vector represented as follows:

$$\Sigma_{\Phi_r} = I_k \otimes \frac{1}{k} [X_p, \hat{x}]^T [X_p, \hat{x}] \quad (28)$$

It is not feasible to find analytic solutions for problems of maximization; however, approximate solutions can be found numerically through optimization [22]. To perform this type of attack, the computational complexity is immense mainly due to $O(m(p+I))$ multiplication of the matrix $X^T X$ along with the Cholesky decomposition.

4.3 Extended MAP attack for random translation

The previous section described the attack on RP which does not take into account extra random translation that is applied in RP method. The application of random translation extends the perturbation model as follows:

$$\overline{[Y_p, y_i]} = \frac{1}{\sqrt{k}\sigma_r} \overline{R[X_p, \hat{x}] + \varphi} \quad (29)$$

Despite this, it is possible to account for random translation by carrying out the translation of both $[X_p, \hat{x}]$ and $[Y_p, y_i]$ such that at least one input-output record pair is aligned at the zero vector or the origin. This pair would represent the result of an RP without translation because the vector itself ($\frac{1}{\sqrt{k}\sigma_r} \overline{R0} = 0$). As the zero vector shows linear dependence on the other vectors, we cannot use that pair for performing an attack on RP. Also, we cannot use the first column of the matrix, which also represents a zero vector. Therefore, this alignment operation can be defined as a function as follows:

$$\alpha \left(\begin{bmatrix} a_{1,1} & a_{1,1} & \dots & a_{1,n} \\ \dots & \dots & \dots & \dots \\ a_{m,1} & a_{m,2} & \dots & a_{m,n} \end{bmatrix} \right) = \begin{bmatrix} a_{1,2} - a_{1,1} & \dots & a_{1,n} - a_{1,1} \\ \dots & \dots & \dots \\ a_{m,2} - a_{m,1} & \dots & a_{m,n} - a_{m,1} \end{bmatrix} \quad (30)$$

We consider that $\alpha(A) = \alpha(A + \varphi)$ for a matrix A with translation φ .

Using α , we can make changes to the equations for MAP attack to account for the random translation as follows:

$$\hat{x}_i = \arg \sup_{\hat{x}} \Phi_r(\alpha(\overline{[Y_p, y_i]})) \quad (31)$$

$$\Sigma_{\Phi_r} = I_k \otimes \frac{1}{k} \alpha[X_p, \hat{x}]^T \alpha[X_p, \hat{x}] \quad (32)$$

This equation accounts for both RP as well as random translation.

4.4 Improved MAP attack(Independent noise)

This section describes the A-RPIN attack, which extends the A-RP by including independent noise as follows:

$$\overline{[Y_p, y_i]} = \frac{1}{\sqrt{k}\sigma_r} \overline{R[X_p, \hat{x}] + \varphi + [\Delta_p, \delta_i]} \quad (33)$$

In this attack, two stages are involved and each stage has its own problem of optimization. In the first stage, the presence of noise ($\alpha(Y_p - \hat{\Delta}_p)$) that have resulted from RP implemented on known input data records as follows:

$$\widehat{\Delta}_p = \arg \sup_{\Delta_p} \frac{1}{kp+1} (\Phi_r(\alpha(\overline{[Y_p - \hat{\Delta}_p]}) + \Sigma \Phi_\delta(\hat{\Delta}_p)) \quad (34)$$

Where, ϕ_δ and $\hat{\Delta}_p \in R^{k,p}$ refer to independent noise which is represented as $(N(0, \sigma_\delta^2))$.

The inclusion of $\frac{1}{kp+1}$ in the equation helps generate resultant attacks like A-RP. Furthermore, \hat{x} has not been included in the equation and so, the covariance matrix for ϕ_r can be represented as follows:

$$\Sigma_{\Phi_r} = I_k \otimes \frac{1}{k} \alpha[X_p]^T \alpha[X_p] \quad (35)$$

In the second stage of this attack, an optimization of the estimated data record \hat{x}_i along with the implemented independent noise $\hat{\delta}_i$ is carried out. In order to carry out this optimization, $\hat{\Delta}_p$ from the first stage is used which balances out the probability of giving rise to the presence of noise ($\alpha(Y_p, y_i) - (\hat{\Delta}_p, \delta)$) are the result of an RP on the input records as follows:

$$\hat{x}_i, \hat{\delta}_i = \arg \sup_{\hat{x}, \hat{\delta}} \frac{1}{2} (\Phi_r(\alpha(Y_p, y_i) - (\hat{\Delta}_p, \delta)) + \frac{1}{k} \Sigma \Phi_\delta(\hat{\delta})) \quad (36)$$

Where, $\hat{x} \in R^m$ l. i. X_p and $\hat{\delta} \in R^k$.

Here, all the mean densities of ϕ_δ are given equal importance. It was seen that when the densities were balanced with attack was enhanced.

When $p = I$, it is considered to be a special case as all $\alpha(X_p)$ is converted into a null matrix. To account for this, the distribution variance $\Phi_\delta(\hat{\delta})$ is increased to two times to $2\sigma_\delta^2$, which takes care of the independent noise in both the known and unknown data records.

4.5 Extended MAP attack for cumulative noise

To account for cumulative noise in the perturbation model, we can modify A-RPIN to A-RPCN as follows:

$$\overline{[Y_p, y_i]} = \frac{1}{\sqrt{k}\sigma_r} \overline{R[X_p, \hat{x}] + \varphi + [\Gamma_p, \gamma_i]} \quad (37)$$

Hence, column of Γ_p is the sum of all columns leading up to the corresponding column of Ω_p .

$$\Omega_p = \begin{matrix} (\omega_{i,j}) \\ k \times p \end{matrix} \quad (38)$$

$$\Gamma_p = \begin{bmatrix} \omega_{1,1} & \sum_{i=1}^{i \leq 2} \omega_{1,i} & \dots & \sum_{i=1}^{i \leq p} \omega_{1,i} \\ \dots & \dots & \dots & \dots \\ \omega_{k,1} & \sum_{i=1}^{i \leq 2} \omega_{k,i} & \dots & \sum_{i=1}^{i \leq p} \omega_{k,i} \end{bmatrix} \quad (39)$$

$$\widehat{\Omega}_p = \arg \sup_{\widehat{\Omega}_p} \frac{1}{kp+1} (\Phi_r(\alpha(\overline{[Y_p - \widehat{\Gamma}_p]}) + \sum \Phi_\omega(\widehat{\Omega}_p)) \quad (40)$$

Where, $\widehat{\Omega}_p \in R^{k,p}$.

The most significant difference lies in using the probability density function based on variations in cumulative noise between different known data records which is represented as follows:

$$\Phi_\omega \sim N(0, (i-h)\sigma_\gamma^2)$$

h denotes the most recent data before the data record i present in X_p .

Both these stream indexes denote the record positions in the complete data stream, and not just in X_p . Let's consider that in X_p , i represents the first column, h is undefined, and the noise is represented by a zero. Therefore, when $p = 1$, no special considerations are necessary as was the case in A-RPIN.

For the second stage, we can derive $\widehat{\Gamma}_p$ from $\widehat{\Omega}_p$ and thereby simultaneously optimize \widehat{x}_i and \widehat{y}_i as follows:

$$\widehat{x}_i, \widehat{y}_i = \arg \sup_{\widehat{x}, \widehat{y}} \frac{1}{2} (\Phi_r(\alpha(\overline{[Y_p, y_i]}) - (\widehat{\Gamma}_p, \widehat{y})) + \frac{1}{k} \sum \Phi_\gamma(\widehat{y})) \quad (41)$$

Where, $\widehat{x} \in R^m$ l. i. X_p and $\widehat{y} \in R^k$ are used for arriving at the probability the formula known as:

$$\Phi_\gamma \sim N(\mu, \sigma^2)$$

$$\mu = \frac{\widehat{y}_h(j-i) + \widehat{y}_j(i-h)}{j-h} \quad (42)$$

$$\sigma^2 = \frac{(i-h)(j-i)\sigma_\gamma^2}{j-h} \quad (43)$$

Where, h denotes index record before i in X_p , and j denotes the stream index of the data record after i in X_p .

If i is the first record of X_p , then:

$$\Phi_\gamma \sim N(\widehat{y}_j, (j-i)\sigma_\gamma^2)$$

Similarly, if i is the last record of X_p , then:

$$\Phi_\gamma \sim N(\widehat{y}_h, (i-h)\sigma_\gamma^2)$$

The matrix operations that result in the complexity of the A-RPIN and A-RPCN are the same, and so, the complexity of computation of the alterations carried out for numerical optimization can be represented as $O(mp+(kp)^3)$.

4.6 Numerical Optimization

It is an unconstrained and non-linear algorithm for numerical optimization. This method requires the optimization of the seeds for all the variables. Therefore, for the variables $\widehat{\delta}$, $\widehat{\omega}$ and \widehat{y} we generated random values based on their respective Gaussian distributions Φ_δ , Φ_ω and Φ_γ . The values of \widehat{x} were initially generated randomly from a uniform distribution which had a range which was the same as the range of each data feature that was present at the

median value (values were input data records of X_p). Our optimization runs for each of the attacks were carried out using three attacks.

While carrying out the optimization runs, a matrix of inputs ($\alpha([X_p, \widehat{x}])$) may be generated wherein all the columns may not show linear independence and the matrix may lack a full column rank. Such combinations of variables were penalized by resulting in the generation of a negative infinity value of the log-probability-density by the objective function.

4.7 Attacks when $p \geq m$

The attacks that have been described above have the requirement of a full column rank of $[X_p, x_i]$; therefore, they can only be used when the number of known data records is less than the features. When the pairs is at least m for a dataset on which RP has been implemented, then any data record x_i must show linear dependence on X_p such that a linear combination of X_p can enable us to recover x_i [8, 9]. However, the relation between ranks of $[Y_p, y_i]$ and $[X_p, x_i]$ is disrupted by noise which might affect the possibility of recovering x_i through a linear combination of X_p . Therefore, when $p \geq m$, more complex attacks need to be considered against RPCN and RPIN.

5. Result and Discussion

We carried out experimental assessment on RP, RPIN, and RPCN. We also evaluated the effects of cumulative noise on privacy and accuracy over a data stream's lifetime. We compared our perturbation methods with other methods.

5.1 Experimental Setup

The efficacy of the methods of privacy-preservation was evaluated by using ϵ -privacy and relative error [23, 24]. The level of success reached in an attempt at record recovery is known as relative error. It represents x_i as follows:

$$\frac{\|x_i - \widehat{x}_i\|}{\|x_i\|}$$

When the privacy of SD is evaluated, known input-output pairs are not taken into consideration as known input-output attacks have not been proposed for SD. Therefore, a naïve attack is used to assess the privacy of SD assuming that the output record that has undergone perturbation is record($\widehat{x}_i = y_i$). As each sensitive value in the data record is subjected to perturbation, it is influenced by only a small set of nearby data records with no participation of the other data records, and so, known input-output attacks will not be effective against SD. In order to carry out an attack, the attacker will need access to the specific set of records that are involved in the perturbation of the target data records which is an unlikely scenario as compared to accessing input-output pairs present in any part of the data stream as in attacks on RP.

The value of σ_r was set to 1 and data dimensionality was not decreased ($k = m$). The impact of random translation was not assessed because all the known input-output attacks are

capable of removing random translation from the perturbation methods.

5.2 Data Set

Table 1: Datasets chosen for our experiment

Dataset	Features	Classes	Records	Real-world?	Stream?	Private?
SEA	3	2	1,00,000	No	Yes	No
RBF	10	5	50,000	No	Yes	No
ELEC	8	2	45,312	Yes	Yes	No
WFR	4	4	5,456	Yes	Yes	No
AREM	6	32	35,999	Yes	Interleaved	Individual's stream
TAXI	7	3	50,000	Yes	Yes	Stream of individuals
POWUSG	10	3	19,735	Yes	Yes	Individual's stream
P2PLNS	10	2	12,682	Yes	Yes	Stream of individuals
PREG	5	2	4,082	Yes	Yes	Stream of individuals
BRCNCR	9	2	10,000	Yes	No	Stream of individuals
ADULT	6	2	32,561	Yes	No	Stream of individuals
HTRU2	8	2	17,898	Yes	No	No

Among our datasets, eight were from the real-world and included ELEC (electricity) [13], WFR eight datasets are considered to be data streams because it is possible to chronologically order their data records; however, we do not have knowledge of the concept drift that might exist in these datasets.

For the datasets POWUSG, TAXI, PREG, and P2PLNS, the sets of data features were decreased to a subset comprising of only numerical features. For POWUSG (amount of power usage) and TAXI (duration of trips), a 3-class classification target was generated using equal-frequency binning for specific target features. PREG (live births versus still-births and miscarriages) and P2PLNS (completed versus defaulted and charged-off loans) were sub-sampled without replacing any feature so that balance was achieved between the classes and classification accuracy could be used as an effective measure of performance. The initial 50,000 data records from TAXI were used for evaluating the efficiency. All the values belonging to the parity feature in PREG were reduced to 1 so that the pregnancy outcome could be erased from the dataset.

The UCI machine learning repository was used to retrieve AREM, BRCNCR, HTRU2, WFR, POWUSG, and ADULT datasets [11].

A lot of these datasets possess sensitive information, which may either be individuals' personal information (P2PLNS,






BRCNCR, TAXI, and PREG) or the data stream may belong to a sensor that monitors individuals (POWUSG and AREM). For instance, in case of TAXI, each data record possesses information about individual taxi trips which includes the time the trip started and checks origin point, which may represent sensitive information such as the home address and/or a personal trip to a health facility (for example). In addition, the information in the dataset may also be confidential for the taxi company as it may provide insights into the operational behavior of the company which may prove hazardous if leaked. Despite this, agency may willingly share their company data so that they can enhance their models for predicting the duration of the trip and assigning prices.

5.3 Comparison of attack type

This will enable us to use the best attack type for comparing the privacy that results from the perturbation methods. Naïve attack is the only attack that is applicable for the SD method and so, an attack type comparison was not performed for this method. In the same way, the ARP attack is the only attack type applicable for the RP method and an attack type comparison was not performed for this method as well. On the other hand, for RPIN and RPCN methods, several attack types are possible. One approach for comparison is to compare attack types separately for RPIN and RPCN. Another approach is to use A-RP in the presence

of additive noise because the expected mean value becomes zero. However, contrary to our expectations, we encountered relative error in the data records that were recovered by RPCN and RPIN which increased with an increase in the pairs. The reason for this may be that in the first stage of both these attacks, the number of variables increases rapidly in the problems of optimization with an increase in $O(kp)$. Hence, variations of the original attack types with just one known input-output pair. When there was more than one known input-output pair, the data record that was nearest to the unknown data record was used.

Table 2: Attack types that were assessed

Perturbation method	Type of attack				
					
RPIN.	AR P	ARP IN	ARPI N-1	MAX(A RP,A- RPIN)	MAX(A RP,A- RPIN-1)
RPCN.	AR P	ARP CN	ARP CN-1	MAX(A RP,A- RPCN)	MAX(A RP,A- RPCN-1)

We compared the attack types by assessing the effectiveness with which the original data records were recovered from the perturbed TAXI dataset (we achieve similar results with the other datasets and therefore, they have been omitted from this discussion). Three perturbed versions of the dataset were generated for each perturbation method, each having varying noise levels. In case of independent noise, we implemented the three different noise levels by using three different values of $\sigma_\delta - 0.05, 0.1, \text{ and } 0.25$. Similarly, for cumulative noise, we adjusted the σ_γ values such that we were able to achieve noise levels depending on records. Furthermore, we used varying numbers of known input-output pairs $- 1, 4\left(\frac{m}{2}\right), \text{ and } 6(m-1)$. For each combination of known input-output pair, noise level, and attack type, we simulated 500 attacks. All attack types had similar performance when only one known input-output pair was used (as can be seen through overlapping points in the graphs). This is because when there is just one data record, then there is no difference in A-RPCN-1 and A-RPIN-1. Also, only attacks that take additive noise into consideration are important in a combined attack because A-RP results in reduced values of probability density. As the effectiveness of the attacks also increases at first, but then decreases when the pairs become maximum.

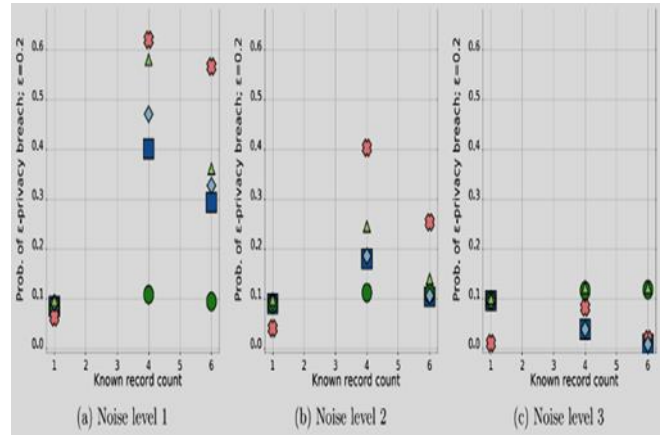


Fig 3: Comparison of different attack types on the RPIN perturbation method at three different noise levels using the TAXI dataset

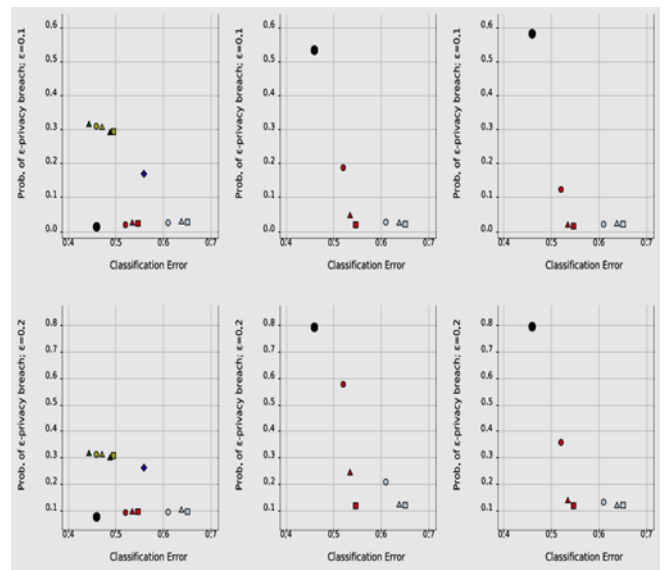


Fig 4: Comparison of different attack types on the RPCN perturbation method at three different noise levels using the TAXI dataset

5.3.1 Comparison of execution time of the attacks

The complexity of computation of the attack types depends on the number of iterations during optimizations and so, the time of execution of the individual attack types were also assessed experimentally. The mean duration of time for 500 attacks across the attack types is shown in figure 5 and the time either increases exponentially or polynomially with an increase in the input output pairs attack types where only one known input-output pair was considered. Hence, even if a greater number of input-output pairs are known to an attacker, it may not enhance the effectiveness of the attack aimed at breaching data privacy where the data is valuable only for a certain duration after it is generated.

5.4 Comparison of Perturbation method

Once the benchmark attacks on privacy were established, we compared the four different perturbation methods across the 12

Datasets. We used three different pairs for RP, RPIN, and RPCN, but not for SD because naïve attacks pair numbers. The values used

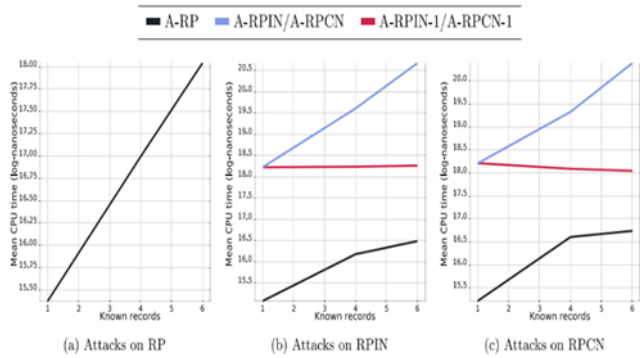


Fig 5: Comparison based on execution time

for pairs were $1, \lfloor \frac{m}{2} \rfloor$, and $m - 1$, with m representing the number of features of a dataset). Three levels of independent and cumulative noise were used for RPIN and RPCN as mentioned above. In case of the SD perturbation method, parameters used for our experiment were the ones reported. The size of the s-window was maintained constant at 3000 data records because only the z-score normalization is affected by this parameter and it does not have a significant effect on privacy. On the other hand, the sizes of the user windows were tested using three different values – 10, 30, and 50. We also varied the SD values to achieve three different levels of noise. SD-30 was tested with 3%, 5%, and 10% SD values, whereas SD-10 and SD-50 were tested with 5% SD value. These variations, however, were not able to result in a level of privacy that was comparable to the other perturbation methods. As a result, we tested the SD method with SD-100 size of user window at 100% SD value representing noise level 4.

The accuracy of the model in learning from the perturbed TAXI dataset is presented in Figure 6, and the legend for the figure is given in Table 3. Perturbation methods that have higher privacy and accuracy are represented by points located nearer to the origin. Therefore, we can see that with an increase in the size of the user window for SD and the levels of noise, there is an improvement in privacy but not in accuracy. These results are significant as they provide insights into the perturbation method’s robustness. Also, when the pairs was one, the privacy was reduced (as seen in Figure 6). The figure also indicates that RPCN performed better than RPIN in terms of the trade-off between accuracy and privacy. For any level of privacy, RPCN had a higher accuracy than RPIN which could be attributed to its gradual implementation of noise unlike RPIN.

We performed analysis so that significant differences between the perturbation methods’ performance could be identified at various levels of noise. These analyses were performed using methods similar to that used for performing comparisons of accuracy and privacy [Privacy-Accuracy

Magnitude (PAM)] was assessed by comparing the probability of a breach of ϵ -privacy and classification error as follows:

$$PAM = \left(\frac{p - p_{min}}{p_{max} - p_{min}} \right)^2 + \left(\frac{e - e_{min}}{e_{max} - e_{min}} \right)^2 \quad (44)$$

Where, p represents $P(\epsilon$ -privacy breach) and e represents error.

Table 3: Legend for figure 6

Noise	RP	RPIN	RPCN	SD-10	SD-30	SD-50	SD-100
None	●						
Level 1		○	●		○		
Level 2		△	▲	▲	▲	▲	
Level 3		□	■		■		
Level 4							◆

Implementation of sum-of-squares favors method and min-max normalization can generate a balance between accuracy and privacy without compromising either of them. We assessed attacks at a ϵ value of 0.2, which is the highest possible value and the most difficult situation for privacy preservation.

Tables 4 and 5 provide the results for privacy and accuracy for all the perturbation methods and the method which resulted in the lowest PAM. Results observed in Figure 6. Among the RPCN variants, RPCN-2 is seen to be the most efficient. Among RPIN variants, RPIN-1 is seen to generate the best accuracy at a reasonable level of privacy. However, for the SD perturbation method is quite sharp. We placed the best perturbation method as we considered privacy to be more important before considering accuracy. Therefore, among the SD variants, we found SD-100-4 to be the most efficient. For the final statistical comparison, we chose RPCN-2, RPIN-1, RP, and SD-100-4 to evaluate differences in the trade-offs between privacy and accuracy among these methods.

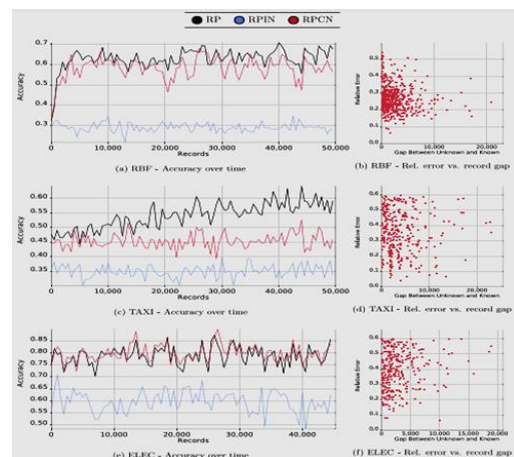


Fig 6: Representation of trade-off among data privacy and accuracy

This hypothesis was rejected at a confidence level of 99%. We then evaluated statistically significant differences among these perturbation methods using the Nemenyi post-hoc analysis. We found a ~ 1.35 critical difference as seen in Figure 7. Our results revealed between privacy was demonstrated by RPCN-2 with a significant improvement over the other methods tested. Therefore, RPCN-2 is the best performing method for striking a balance between accuracy and privacy.

Table 4: Classification errors for the various perturbation methods on five different datasets.

Method	TA XI.	POW USG.	P2PL NS.	PRE G.	BRCN CR.
RP	0.459	0.310	0.350	0.381	0.014
RPIN-1	0.610	0.452	0.361	0.407	0.026
RPIN-2	0.638	0.498	0.377	0.439	0.041
RPIN-3	0.651	0.552	0.411	0.471	0.106
RPCN-1	0.520	0.307	0.350	0.380	0.018
RPCN-2	0.534	0.306	0.353	0.389	0.024
RPCN-3	0.546	0.310	0.364	0.413	0.032
SD-10-2	0.443	0.309	0.374	0.433	0.035
SD-30-1	0.459	0.313	0.375	0.449	0.034
SD-30-2	0.471	0.314	0.381	0.453	0.034
SD-30-3	0.495	0.317	0.379	0.461	0.035
SD-50-2	0.488	0.320	0.379	0.451	0.036
SD-100-4	0.559	0.360	0.395	0.466	0.037

The numbers in the perturbation methods indicate the noise levels. The numbers in bold indicate the minimum PAM value that was generated by a method for that particular dataset.

Table 5: Probability of breach of ϵ -privacy ($\epsilon = 0.2$, input output pairs = $m - 1$) for the various perturbation methods on five different datasets.

Method	TA XI.	POW USG.	P2PL NS.	PRE G.	BRCN CR.
RP	0.796	0.738	0.840	0.322	0.306
RPIN-1	0.132	0.254	0.112	0.052	0.006
RPIN-2	0.118	0.224	0.112	0.016	0.002
RPIN-3	0.120	0.224	0.108	0.016	0.006
RPCN-1	0.358	0.566	0.196	0.180	0.074
RPCN-2	0.136	0.352	0.112	0.076	0.014
RPCN-3	0.118	0.210	0.096	0.018	0.004
SD-10-2	0.314	0.574	0.398	0.302	0.414
SD-30-1	0.312	0.574	0.388	0.280	0.332
SD-30-2	0.310	0.574	0.370	0.246	0.254
SD-30-3	0.308	0.574	0.344	0.162	0.154
SD-50-2	0.298	0.574	0.358	0.180	0.176
SD-100-4	0.262	0.562	0.184	0.018	0.002

The numbers in bold indicate the minimum PAM value that was generated by a method for that particular dataset.

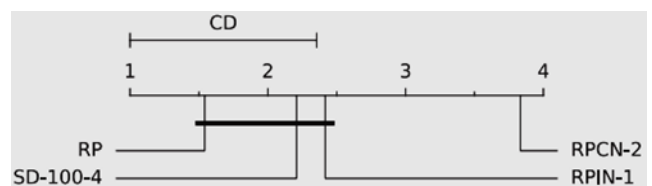


Fig 7: Critical difference for the evaluation of tradeoff between different methods

5.5 Analysis of trends for perturbation of cumulative noise

For the RPCN method, the cumulative noise keeps increasing and therefore, privacy and accuracy may be affected over time. These trends are given in Figure 8 for TAXI, RBF, and ELEC datasets. As the trends were similar across all datasets, they have been eliminated from this discussion.

Figure 8: Trends in privacy and accuracy over time for datasets perturbed using the RPCN method

This indicates that, over time, the value of input output will decrease as the privacy level of RPCN increases. It should be noted that when cumulative noise is present, the accuracy remains stable over time and it does not decrease even when the noise level increases. This trend was specifically seen in the case of the RBF dataset which demonstrates continuous drift. Hence, according to these results, the ARF classifier can not only adapt to increasing noise levels but also to concept drift.

Figure 9 demonstrates how accuracy may be maintained by plotting a line which determines each tree's depth in the ARF ensemble over time for the TAXI dataset perturbed using the various methods based on RP at a noise level of 3. The mean depth of the tree was calculated for every 100 data records.

Therefore, the analysis of trends demonstrated that the presence of cumulative noise could improve privacy over time and the gradual addition of noise could be considered as concept drift by the ARF classifier to maintain accuracy at a stable level.

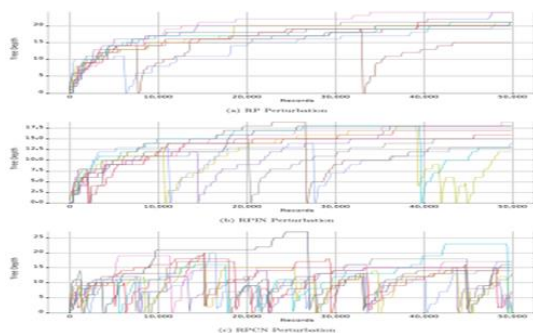


Fig 9: Changes in trends of depth of tree generated by the ARF classifier for the TAXI dataset perturbed by RPCN, RPIN, and RP

6. Conclusion and Future Work

In this work we used combined techniques such as random translation and random projection (RP), and independent noise (for RPIN) for (RPCN). We also developed variations of the MAP attacks that we implemented against the perturbation methods. We also showed that the best attack against the RPCN method was a combination of two attacks— one which accounted for cumulative noise and the other which did not account for cumulative noise. This attack was not as effective for records farther away from the known data records as compared to records that were nearer to the known data records indicating that with the RPCN method, the privacy of data gradually increased over time. Our findings have a lot of scope for future research. We have concentrated on tasks related to classification; however, the perturbation methods can also be used for other tasks such as clustering, detection of anomalies, and regression. Our

proposed method can be used for numerical data by considering them to be integer; however, if someone knows that the data is numerical, then the privacy becomes compromised [1, 3]. Therefore, our method needs to be improved to ensure a higher level of privacy for nominal data.

References

- [1] Bifet A, Kirkby R. Data stream mining: a practical approach. The university of Waikato. 2009 Aug. Centers for Disease Control and Prevention. National survey of family growth data. 2005 (retrieved February 12, 2019) <http://www.greenteapress.com/thinkstats/nsfg.html>.
- [2] Chamikara MAP, Bertók P, Liu D, Camtepe S, Khalil I. An efficient and scalable privacy preserving algorithm for big data and data streams. *Comput Secur* 2019;87:101570. <http://dx.doi.org/10.1016/j.cose.2019.101570>.
- [3] Chamikara MAP, Bertók P, Liu D, Camtepe S, Khalil I. An efficient and scalable privacy preserving algorithm for big data and data streams. *Comput Secur* 2019;87:101570. <http://dx.doi.org/10.1016/j.cose.2019.101570>.
- [4] Matatov N, Rokach L, Maimon O. Privacy-preserving data mining: A feature set partitioning approach. *Information Sciences*. 2010 Jul 15;180(14):2696-720. <https://doi.org/10.1016/j.ins.2010.03.011>
- [5] Chamikara MAP, Bertók P, Liu D, Camtepe S, Khalil I. Efficient privacy preservation of big data for accurate data mining. *Inform Sci* 2020;527:420–43. <http://dx.doi.org/10.1016/j.ins.2019.05.053>.
- [6] Chamikara MAP, Bertok P, Khalil I, Liu D, Camtepe S. Privacy preserving distributed machine learning with federated learning. *Comput Commun* 2021;171:112–25. <http://dx.doi.org/10.1016/j.comcom.2021.02.014>
- [7] Chamikara MA, Bertok P, Liu D, Camtepe S, Khalil I. Efficient privacy preservation of big data for accurate data mining. *Information Sciences*. 2020 Jul 1;527:42043. <https://doi.org/10.1016/j.ins.2019.05.053>
- [8] Virupaksha S, Dondeti V. Anonymized noise addition in subspaces for privacy preserved data mining in high dimensional continuous data. *Peer-to-Peer Networking and Applications*. 2021 May;14(3):1608-28. <https://doi.org/10.1007/s12083-021-01080-y>
- [9] Deshkar PA, Patil JM, Niranjane PB, Niranjane V, Thakur N, Dabhade VD. Studies on the Use of Various Noise Strategies for Perturbing Data in Privacy-Preserving Data Mining. *International Journal of*

- [10] Denham B, Pears R, Naeem MA. Enhancing random projection with independent and cumulative additive noise for privacy-preserving data stream mining. *Expert Syst Appl* 2020;152(8):321–35. <http://dx.doi.org/10.1016/j.eswa.2020.113380>.
- [11] Virupaksha S, Dondeti V. Subspace based noise addition for privacy preserved data mining on high dimensional continuous data. *Journal of Ambient Intelligence and Humanized Computing*. 2020 Mar 21;1-7. <https://doi.org/10.1007/s10618-005-1396-1>
- [12] Fang W, Wen XZ, Zheng Y, Zhou M. A survey of big data security and privacy preserving. *IETE Tech Rev* 2017; 34(5):544–60. <http://dx.doi.org/10.1080/02564602.2016.1215269>.
- [13] Kadampur MA. A noise addition scheme in decision tree for privacy preserving data mining. *arXiv preprint arXiv:1001.3504*. 2010 Jan 20. <https://doi.org/10.48550/arXiv.1001.3504>
- [14] K. Xing, C. Hu, J. Yu, X. Cheng and F. Zhang, "Mutual Privacy Preserving k -Means Clustering in Social Participatory Sensing," in *IEEE Transactions on Industrial Informatics*, vol. 13, no. 4, pp. 2066-2076, Aug. 2017, doi: 10.1109/TII.2017.2695487.
- [15] J. Vaidya, B. Shafiq, W. Fan, D. Mehmood and D. Lorenzi, "A Random Decision Tree Framework for Privacy-Preserving Data Mining," in *IEEE Transactions on Dependable and Secure Computing*, vol. 11, no. 5, pp. 399-411, Sept.-Oct. 2014, doi: 10.1109/TDSC.2013.43.
- [16] Z. Xiao, X. Fu and R. S. M. Goh, "Data Privacy-Preserving Automation Architecture for Industrial Data Exchange in Smart Cities," in *IEEE Transactions on Industrial Informatics*, vol. 14, no. 6, pp. 2780-2791, June 2018, doi: 10.1109/TII.2017.2772826.
- [17] H. Chen, K. Mei, Y. Zhou, N. Wang, M. Tang and G. Cai, "A Density Peaking Clustering Algorithm for Differential Privacy Preservation," in *IEEE Access*, vol. 11, pp. 54240-54253, 2023, doi: 10.1109/ACCESS.2023.3281652.
- [18] T. Tassa and D. J. Cohen, "Anonymization of Centralized and Distributed Social Networks by Sequential Clustering," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 25, no. 2, pp. 311-324, Feb. 2013, doi: 10.1109/TKDE.2011.232.
- [19] M. Kanmaz, M. A. Aydin and A. Sertbas, "A New Geometric Data Perturbation Method for Data Anonymization Based on Random Number Generators," in *Journal of Web Engineering*, vol. 20, no. 6, pp. 1947-1970, September 2021, doi: 10.13052/jwe1540-9589.20613.
- [20] K. Bhaduri, M. D. Stefanski and A. N. Srivastava, "Privacy-Preserving Outlier Detection Through Random Nonlinear Data Distortion," in *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 41, no. 1, pp. 260-272, Feb. 2011, doi: 10.1109/TSMCB.2010.2051540.
- [21] Y. -T. Tsou, H. -L. Chen and J. -Y. Chen, "RoD: Evaluating the Risk of Data Disclosure Using Noise Estimation for Differential Privacy," in *IEEE Transactions on Big Data*, vol. 7, no. 1, pp. 214-226, 1 March 2021, doi: 10.1109/TBDATA.2019.2916108.
- [22] K. -P. Lin and M. -S. Chen, "On the Design and Analysis of the Privacy-Preserving SVM Classifier," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 23, no. 11, pp. 1704-1717, Nov. 2011, doi: 10.1109/TKDE.2010.193.
- [23] S. M. Darwish, R. M. Essa, M. A. Osman and A. A. Ismail, "Privacy Preserving Data Mining Framework for Negative Association Rules: An Application to Healthcare Informatics," in *IEEE Access*, vol. 10, pp. 76268-76280, 2022, doi: 10.1109/ACCESS.2022.3192447.
- [24] L. Li, R. Lu, K. -K. R. Choo, A. Datta and J. Shao, "Privacy-Preserving-Outsourced Association Rule Mining on Vertically Partitioned Databases," in *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 8, pp. 1847-1861, Aug. 2016, doi: 10.1109/TIFS.2016.2561241.
- [25] R. Mendes and J. P. Vilela, "Privacy-Preserving Data Mining: Methods, Metrics, and Applications," in *IEEE Access*, vol. 5, pp. 10562-10582, 2017, doi: 10.1109/ACCESS.2017.2706947.
- [26] Z. Zhou, Y. Wang, X. Yu and J. Miao, "A Targeted Privacy-Preserving Data Publishing Method Based on Bayesian Network," in *IEEE Access*, vol. 10, pp. 89555-89567, 2022, doi: 10.1109/ACCESS.2022.3201641.
- [27] Y. Li, M. Chen, Q. Li and W. Zhang, "Enabling Multilevel Trust in Privacy Preserving Data Mining," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 9, pp. 1598-1612, Sept. 2012, doi: 10.1109/TKDE.2011.124.
- [28] M. B. Malik, M. A. Ghazi and R. Ali, "Privacy Preserving Data Mining Techniques: Current Scenario and Future Prospects," 2012 Third International Conference on Computer and Communication Technology, Allahabad, India, 2012, pp. 26-32, doi: 10.1109/ICCCT.2012.15.