

Utility Preserving Secure Private Data Release

Jasjeet Dhaliwal¹ Geoffrey So¹ Aleatha Parker-Wood¹ Melanie Beck¹

Abstract

Differential privacy mechanisms that also make reconstruction of the data impossible come at a cost - a decrease in utility. In this paper, we tackle this problem by designing a private data release mechanism that makes reconstruction of the original data impossible and also preserves utility for a wide range of machine learning algorithms. We do so by combining the Johnson-Lindenstrauss (JL) transform with noise generated from a Laplace distribution. While the JL transform can itself provide privacy guarantees (Blocki et al., 2012) and make reconstruction impossible, we do not rely on its differential privacy properties and only utilize its ability to make reconstruction impossible. We present novel proofs to show that our mechanism is differentially private under single element changes as well as single row changes to any database. In order to show utility, we prove that our mechanism maintains pairwise distances between points in expectation and also show that its variance is proportional to the dimensionality of the subspace we project the data into. Finally, we experimentally show the utility of our mechanism by deploying it on the task of clustering.

1. Introduction

While the recent surge in data available for machine learning has created new opportunities, it also poses a risk to the privacy of individuals. Differential privacy (Dwork et al., 2006) is the most widely accepted framework that attempts to address this issue by capturing precisely how much additional information of an individual is leaked by participating in a database that would not have been leaked otherwise. Consequently, by providing guarantees within the framework of differential privacy, one can protect the privacy of the individuals participating in a database while utilizing

the private data to build machine learning models.

However, differentially private releases of aggregate statistics have been shown to be susceptible to reconstruction and tracing attacks (Dwork et al., 2017) under certain constraints. This is also true in the non-interactive setting where the amount of noise added to the data must be increased drastically in order to make reconstruction difficult (Bhowmick et al., 2018; Dwork et al., 2014). Even differentially private machine learning models have been shown to be susceptible to membership inference attacks that leak information about individual participation (Rahman et al.; Shokri et al., 2017). The above issues create the need for stronger privacy mechanisms that make it provably impossible for an adversary to reconstruct the original data, hence eliminating concerns about privacy. However, mechanisms that do so, achieve this goal at the cost of a significant drop in utility (Blocki et al., 2012). One approach to improve this utility has been to design private data release mechanisms that are tailored to specific machine learning algorithms (Upadhyay, 2014; Blocki et al., 2012; Sheffet, 2015b). But this approach limits the ability of the analyst to compare different machine learning algorithms on a given dataset. For instance, if an analyst wants to perform clustering in order to better understand the data before classification, the covariance matrix or other similar aggregate statistics may not suffice. Similarly, an analyst may wish to compare a neural network, linear regression, and a random forest on a given dataset before deploying a model. It is therefore important to design a mechanism that privately releases data in a form that it can be used by a wide variety of machine learning algorithms while preserving utility.

In this paper, we tackle this problem of preserving utility for a broad class of machine learning algorithms while also making reconstruction of the original data impossible. We do so by using the JL transform in conjunction with the Laplace mechanism. The JL transform is a powerful dimensionality reduction tool that preserves pairwise distances between points (Dasgupta & Gupta, 2003) and has also been shown to provide differential privacy guarantees by itself (Blocki et al., 2012). We do not rely on the differential privacy guarantees of the JL transform but utilize the fact that it makes it impossible to reconstruct the exact original values (or the original dimensionality) of the data (Liu et al., 2006). We do not use the differential privacy

¹Center for Advanced Machine Learning, Symantec Corporation, Mountain View, California, USA. Correspondence to: Jasjeet Dhaliwal <jasjeet.dhaliwal@symantec.com>.

properties of the JL transform because in order to achieve differential privacy via the JL transform by itself, certain constraints must be imposed on the rank and the spectrum of the data matrix. More specifically, the data matrix is required to be full rank and the smallest eigenvalue of the data matrix must be larger than a given threshold. For data matrices that do not meet the required constraints, the mechanism in (Blocki et al., 2012) modifies the spectrum of the matrix via an operation that greatly compromises utility. This leaves open the problem of utilizing the ability of the JL transform to make reconstruction impossible while providing differential privacy via a mechanism that preserves utility.

We make progress in this direction by combining the power of the JL transform and the Laplace mechanism in a manner that preserves utility while also providing differential privacy guarantees such that the original data cannot be reconstructed. Our approach is similar to (Kenthapadi et al., 2012) and can be considered as an extension to their work. Our data release mechanism provides utility for general machine learning tasks and is differentially private under single element changes as well as row changes. We prove the differential privacy guarantees provided by our mechanism and also propose an algorithm that maintains pairwise distances between private data points in expectation. We chose the most general task of clustering in order to show the effectiveness of our methods experimentally.

Our contributions in this paper are:

- We propose a differentially private data release mechanism that preserves privacy and makes it impossible for an adversary to reconstruct the original data.
- We propose a distance recovery algorithm and prove that it maintains pairwise distances in expectation. We also prove precise variance guarantees for the distance recovery algorithm.
- We experimentally validate the utility of our mechanism by showing that it maintains pairwise distances and performs well on the general task of clustering

The paper is organized as follows: we first provide the required background on differential privacy in Section 2. We describe our mechanism in Section 3 and prove its differential privacy guarantees. We provide utility guarantees in Section 4 and experimentally validate our claims in Section 5, showing the effectiveness of our mechanism. Related work is covered in Section 6, followed by the Conclusion in Section 7.

2. Background

In this section we define differential privacy and also cover other necessary mathematical background required for our results.

2.1. Differential Privacy

Differential privacy captures precisely how likely is it for a third-party to ascertain whether an individual participated in a database or not. In order to formalize the definition of differential privacy, we first introduce the notion of neighboring databases.

Definition 1 (Neighboring databases). Given an input space $\mathcal{X} \subseteq \mathbb{R}^d$, we can represent a database with n entries, $X \in \mathcal{X}^n$ as $X \in \mathbb{R}^{n \times d}$. Then, two databases $X_1, X_2 \in \mathbb{R}^{n \times d}$ are row-wise neighbors if they differ in exactly one row. They are considered element-wise neighbors, if they differ in exactly one element.

Definition 2 (Probability Simplex). Given a set \mathcal{Y} , the probability simplex over \mathcal{Y} is defined as : $\Delta\mathcal{Y} = \{y \in \mathbb{R}^{|\mathcal{Y}|} : y_i \geq 0, \sum_{i=1}^{|\mathcal{Y}|} y_i = 1\}$

Definition 3 (Randomization Mechanism). Given two sets \mathcal{X}, \mathcal{Y} , a randomization mechanism is a function $\mathcal{M} : \mathcal{X} \rightarrow \Delta\mathcal{Y}$.

Thus, a randomization mechanism defines a probability distribution over the set \mathcal{Y} . Given an input $x \in \mathcal{X}$, a randomization mechanism \mathcal{M} , maps x to $y \in \mathcal{Y}$ with probability $(\mathcal{M}(x))_y$, which is the probability for element y under the distribution $(\mathcal{M}(x))$.

Definition 4 (Privacy Loss). For a randomization mechanism \mathcal{M} , the privacy loss for two neighboring databases X_1, X_2 , is defined as: $\mathcal{L}_{\mathcal{M}(X_1)||\mathcal{M}(X_2)}^D = \ln \left(\frac{\mathbb{P}[\mathcal{M}(X_1) \subseteq D] - \delta}{\mathbb{P}[\mathcal{M}(X_2) \subseteq D]} \right)$

Definition 5 (Differential Privacy). For any $\epsilon > 0$, and $\delta \in [0, 1]$, a randomization mechanism \mathcal{M} is (ϵ, δ) differentially private on domain \mathcal{X} if for two neighboring databases X_1, X_2 , the privacy loss $\left| \mathcal{L}_{\mathcal{M}(X_1)||\mathcal{M}(X_2)}^D \right| \leq \epsilon$.

For a more thorough review of (ϵ, δ) privacy, the reader is referred to (Dwork et al., 2014).

2.2. Johnson-Lindenstrauss Lemma

The Johnson-Lindenstrauss Lemma states that a set of points in a high-dimensional space can be embedded into a lower dimensional space such that the distances between

the projected points are preserved with high probability. We provide a statement of this lemma that relies on a projection matrix using values from the Gaussian distribution.

Consider a finite set $S \subset \mathbb{R}^d$ with $|S| = n$. Let $P \in \mathbb{R}^{d \times k}$ be a real valued matrix such that $P_{ij} \sim \mathcal{N}(0, \frac{1}{k})$, where $k = \Omega(\Lambda^{-2} \log(n))$ for $0 < \Lambda \leq 1$. Then for any $x, y \in S$, we have:

$$(1 - \Lambda) \|x - y\|_2^2 \leq \|xP - yP\|_2^2 \leq (1 + \Lambda) \|x - y\|_2^2$$

Further, $\mathbb{E} [\|xP - yP\|_2^2] = \|x - y\|_2^2$. This result is called the Johnson-Lindenstrauss (JL) lemma. We refer the reader to (Dasgupta & Gupta, 2003) for a proof of the lemma.

3. Privacy Guarantees

Our randomization mechanism utilizes the JL transform to reduce the dimensionality of the input and then uses the Laplacian mechanism to provide differential privacy guarantees. Since the elements of the JL matrix are normally distributed, we utilize their properties in conjunction with the Laplace mechanism to provide differential privacy guarantees while still maintaining utility. We do not rely on the differential privacy properties of the JL transform itself (Blocki et al., 2012) because it requires a transformation of the data matrix that does not preserve any utility in certain practical cases (see Section 6 for details).

Instead our mechanism design follows that of (Kenthapadi et al., 2012) with two key differences: a) our mechanism adds noise from the Laplace distribution (whereas (Kenthapadi et al., 2012) added noise from a Gaussian distribution) b) our mechanism provides privacy guarantees with respect to element and row-wise changes (whereas (Kenthapadi et al., 2012) only provide guarantees with respect to element wise changes). The JL transform not only reduces dimensionality of the input, but also provides further security from attackers by making it impossible to reconstruct the original data values if the JL transformation matrix is kept secret (Liu et al., 2006). We now describe our randomization mechanism.

3.1. Randomization Mechanism

Given database $X \in \mathbb{R}^{n \times d}$, our mechanism first projects the data onto a lower dimensional subspace \mathbb{R}^k , with $k \ll d$, and then adds a noise matrix $\Delta \in \mathbb{R}^{n \times k}$ to the projected data. The entries of this noise matrix are drawn i.i.d from a Laplacian distribution. The mechanism requires the projection parameter k which determines the dimensionality of the subspace that we wish to project the data into. In addition, it requires the privacy parameters c and ϵ in order to determine the scale of the Laplacian distribution, where ϵ is determined by the level of privacy we wish to maintain and c is a parameter that will become

clear in the proofs of privacy guarantees. Algorithm 1. outlines our mechanism.

Algorithm 1: Randomization Mechanism

Input : $X \in \mathbb{R}^{n \times d}, k, c, \epsilon$

Output : $Z \in \mathbb{R}^{n \times k}$

1. Construct JL projection matrix $P \in \mathbb{R}^{d \times k}$ such that $P_{ij} \sim \mathcal{N}(0, \frac{1}{k})$
 2. Set $Y = XP$
 3. Construct noise matrix $\Delta \in \mathbb{R}^{n \times k}$ such that $\Delta_{ij} \sim \text{Laplacian}(0, \frac{c}{\epsilon})$.
 4. Return $Z = Y + \Delta$
-

Note that we do not release the projection matrix P , in order to eliminate the possibility of a reconstruction attack (Dwork & Yekhanin, 2008).

3.2. Privacy Guarantees

Lemma 1. For any $X, X' \in \mathbb{R}^{n \times d}$, such that X and X' differ in exactly one element with $\|X - X'\|_1 \leq 1$, we have for any $A \in \mathbb{R}^{d \times k}$, $\|XA - X'A\|_1 \leq \sqrt{k} \max_{1 \leq i \leq d} \|A_i\|_2$, where A_i is the i th row of A .

Proof. We prove the above by direct calculation.

$$\begin{aligned} \|XA - X'A\|_1 &= \|(X - X')A\|_1 \\ &\leq \max_{1 \leq i \leq d} \sum_{j=1}^k |A_{ij}| \\ &= \max_{1 \leq i \leq d} \|A_i\|_1 \\ &\leq \sqrt{k} \max_{1 \leq i \leq d} \|A_i\|_2 \end{aligned}$$

where the second inequality follows from the fact that $(X - X')$ only contains one non-zero element and the last inequality follows from the Cauchy-Schwarz inequality for inner product spaces. \square

Lemma 2. (Kenthapadi et al., 2012) Let $P \in \mathbb{R}^{d \times k}$ such that $P_{ij} \sim \mathcal{N}(0, \frac{1}{k})$, then $\Pr \left[\max_{1 \leq i \leq d} \|P_i\|_2 > 1 + \sqrt{\frac{2x}{k}} \right] < de^{-x}$, for any $x > 0$.

Theorem 3. For any two element-wise neighboring databases $X, X' \in \mathbb{R}^{n \times d}$, such that $\|X - X'\|_1 \leq 1$, Algorithm 1. achieves ϵ -differential privacy with probability at least $1 - de^{-\frac{k}{2}}$, with respect to changes in a single element.

Proof. Using Lemma 1 we have, $\|XP - X'P\|_1 \leq \sqrt{k} \max_{1 \leq i \leq d} \|P_i\|_2$. Let, $Y = XP$ and $Y' = X'P$ and without loss of generality, consider $Y, Y', \Delta \in \mathbb{R}^{nk}$. Now, set $Z = Y + \Delta$, $Z' = Y' + \Delta$, and let $D \subset \mathbb{R}^{nk}$. Due to the i.i.d assumption on the elements of Δ we have:

$$\begin{aligned} \Pr[Z \in D] &= \frac{1}{(2b)^{nk}} \int_D e^{-\frac{1}{b}(\|z-Y\|_1)} dz \\ &\geq \frac{1}{(2b)^{nk}} \int_D e^{-\frac{1}{b}(\|z-Y'\|_1 + \|Y'-Y\|_1)} dz \\ &= \frac{1}{(2b)^{nk}} \int_D e^{-\frac{1}{b}(\|z-Y'\|_1)} e^{-\frac{1}{b}(\|Y'-Y\|_1)} dz \\ &= e^{-\frac{1}{b}(\|Y'-Y\|_1)} \Pr[Z' \in D] \end{aligned}$$

$$\implies \Pr[Z' \in D] \leq e^{\frac{1}{b}(\|Y'-Y\|_1)} \Pr[Z \in D]$$

We set $\epsilon = \frac{c}{b}$ in Algorithm 1, therefore, in order to preserve privacy, we must constrain $\frac{\|Y-Y'\|_1}{b} < \frac{c}{b}$.

$$\begin{aligned} \Pr\left[\frac{\|Y-Y'\|_1}{b} > \frac{c}{b}\right] &= \Pr[\|Y-Y'\|_1 > c] \\ &\leq \Pr\left[\max_{1 \leq i \leq d} \|P_i\|_2 > \frac{c}{\sqrt{k}}\right] \end{aligned}$$

Setting $c = 2\sqrt{k}$, and x in Lemma 2 to $\frac{k}{2}$, we get:

$$\Pr\left[\max_{1 \leq i \leq d} \|P_i\|_2 > 2\right] \leq de^{-\frac{k}{2}}$$

Hence, Algorithm 1. achieves ϵ -differential privacy with probability at least $1 - de^{-\frac{k}{2}}$. \square

Thus, Theorem 3 provides privacy guarantees and also utilizes the dimensionality reduction properties of the JL transform. It is similar to the work done by (Kenthapadi et al., 2012) in which the authors define a random mechanism that first does a JL transform and then adds Gaussian noise. Since the above result is limited to providing privacy guarantees for element-wise changes, we now focus on extending it to row-wise changes. That is, we now show that Algorithm 1. provides differential privacy when two databases differ by one row.

Lemma 4. If $P \in \mathbb{R}^{d \times k}$ with $P_{ij} \sim \mathcal{N}(0, \frac{1}{k})$, and $v \in \mathbb{R}^d$ then $\Pr\left[k \max_{1 \leq i \leq k} \left|\sum_{j=1}^d v_j P_{ji}\right| > kt\right] \leq 2ke^{\frac{-kt^2}{2\|v\|_2^2}}$.

Proof. First note that, $\sum_{j=1}^d v_j P_{ji} \sim \mathcal{N}(0, \frac{\|v\|_2^2}{k})$. Therefore, it is Gaussian and hence Sub-Gaussian, which lets us use tail bounds for Sub-Gaussian random variables and get:

$$\Pr\left[\left|\sum_{j=1}^d v_j P_{ji}\right| > t\right] \leq 2e^{\frac{-kt^2}{2\|v\|_2^2}}. \text{ Using the union bound and multiplying both sides by } k, \text{ we get the desired result. } \square$$

We now show that our mechanism provides differential privacy guarantees with respect to row changes in the data matrix.

Theorem 5. For any two row-wise neighboring databases $X, X' \in \mathbb{R}^{n \times d}$, that differ in row m , such that $\|X_m - X'_m\|_2^2 \leq \alpha$, Algorithm 1. achieves ϵ -differential privacy with probability at least $1 - 2ke^{\frac{-kt^2}{2\alpha}}$, where $t \geq \sqrt{\frac{2\ln 2k}{k}}\alpha$.

Proof. Suppose that X, X' differ in row m , then we have:

$$\begin{aligned} \|XP - X'P\|_1 &= \sum_{i=1}^k \left| \sum_{j=1}^d P_{ji}(X_{mj} - X'_{mj}) \right| \\ &\leq k \max_{1 \leq i \leq k} \left| \sum_{j=1}^d P_{ji}(X_{mj} - X'_{mj}) \right| \end{aligned}$$

Next, we can see that $\sum_{j=1}^d P_{ji}(X_{mj} - X'_{mj}) \sim \mathcal{N}(0, \frac{\|X_m - X'_m\|_2^2}{k})$. Hence, we can use Lemma 4 to get:

$$\begin{aligned} \Pr\left[k \max_{1 \leq i \leq k} \left| \sum_{j=1}^d P_{ji}(X_{mj} - X'_{mj}) \right| > kt\right] \\ \leq 2ke^{\frac{-kt^2}{2\|X_m - X'_m\|_2^2}} \end{aligned}$$

Setting $c = kt$, and $t \geq \sqrt{\frac{2\ln 2k}{k}}\alpha$ we ensure that

$2ke^{\frac{-kt^2}{2\|X_m - X'_m\|_2^2}} \in [0, 1]$. Once again letting $Y = XP$ and $Y' = X'P$ and without loss of generality, letting $Y, Y', \Delta \in \mathbb{R}^{nk}$, we set $Z = Y + \Delta$, $Z' = Y' + \Delta$, and let $D \subset \mathbb{R}^{nk}$. By following the same steps as we did in Theorem 3. we get:

$$\Pr[Z' \in D] \leq e^{\frac{1}{b}(\|Y'-Y\|_1)} \Pr[Z \in D]$$

Now, we constrain $\frac{\|Y-Y'\|_1}{b} < \frac{c}{b}$.

$$\begin{aligned}
 & \Pr \left[\frac{\|Y - Y'\|_1}{b} > \frac{c}{b} \right] \\
 &= \Pr [\|Y - Y'\|_1 > c] \\
 &= \Pr [\|(X - X')P\|_1 > c] \\
 &\leq \Pr \left[k \max_{1 \leq i \leq k} \left| \sum_{j=1}^d P_{ji}(X_{mj} - X'_{mj}) \right| > c \right] \\
 &\leq 2ke^{\frac{-kt^2}{2\|X_m - X'_m\|_2^2}} \\
 &\leq 2ke^{\frac{-kt^2}{2\alpha}}
 \end{aligned}$$

Hence, Algorithm 1. achieves ϵ -differential privacy with probability at least $1 - 2ke^{\frac{-kt^2}{2\alpha}}$.

□

4. Utility Guarantees

A differentially private mechanism that is also an isometric isomorphism would allow any machine learning algorithm to extract the same amount of utility from the private data as it could from the non-private data. Drawing from that intuition, we also measure utility as was proposed in (Kenthapadi et al., 2012), by the degree to which a privacy mechanism preserves pairwise distances after its action. That is, a mechanism that allows pairwise distances to be preserved in the private representation of the data is more useful than one that does not. In order to capture this notion, we first define a distance recovery algorithm in Algorithm 2 that takes as input, two private data points and outputs the distance between them. We then show that the algorithm preserves squared distances in expectation. Further, we show that the variance of the squared distance between any two points is proportional to the dimensionality of the subspace that the mechanism projects the data into.

Algorithm 2: Recover Distance

Input : $Z \in \mathbb{R}^{n \times k}$, σ^2 , $(i, j) \in \{1, \dots, n\} \times \{1, \dots, n\}$

Output : Distance between Z_i and Z_j

1. Output $\mathcal{D}(Z_i, Z_j) = \|Z_i - Z_j\|_2^2 - 2k\sigma^2$
-

4.1. Guarantees

Claim 6. Let $S \subset \mathbb{R}^d$ with $|S| = n$. Then, given any two entries in this set $x_i, x_j \in S$, let $y_i = x_iP + \Delta_i$ and $y_j = x_jP + \Delta_j$, where P and Δ_i, Δ_j are the projection matrix and the noise vectors respectively. Let, $\mathcal{D}(\cdot, \cdot)$ be defined as in Algorithm 2. Then, $\mathcal{D}(y_i, y_j)$ is an unbiased estimator of $\|x_i - x_j\|_2^2$.

Proof. Let $\Delta = \Delta_i - \Delta_j$, and let σ^2 be the variance of the entries of the projection matrices Δ_i and Δ_j . Then we have:

$$\begin{aligned}
 \mathbb{E}[\mathcal{D}(y_i, y_j)] &= \mathbb{E}[\|x_iP + \Delta_i - x_jP - \Delta_j\|_2^2 - 2k\sigma^2] \\
 &= \mathbb{E}[\|(x_i - x_j)P + \Delta\|_2^2 - 2k\sigma^2] \\
 &= \mathbb{E}[\|(x_i - x_j)P\|_2^2 + \|\Delta\|_2^2 + 2\langle (x_i - x_j)P, \Delta \rangle - 2k\sigma^2] \\
 &= \mathbb{E}[\|(x_i - x_j)P\|_2^2] + \mathbb{E}[\|\Delta\|_2^2] + \\
 &\quad 2\mathbb{E}[\langle (x_i - x_j)P, \Delta \rangle] - 2k\sigma^2 \\
 &= \mathbb{E}[\|(x_i - x_j)P\|_2^2] + \mathbb{E}\left[\sum_{t=1}^k (\Delta_{it} - \Delta_{jt})^2\right] + \\
 &\quad 2\mathbb{E}[\langle (x_i - x_j)P, \Delta \rangle] - 2k\sigma^2 \\
 &= \mathbb{E}[\|(x_i - x_j)P\|_2^2] + \mathbb{E}\left[\sum_{t=1}^k (\Delta_{it})^2 + (\Delta_{jt})^2 - 2\Delta_{it}\Delta_{jt}\right] + \\
 &\quad 2\mathbb{E}[\langle (x_i - x_j)P, \Delta \rangle] - 2k\sigma^2 \\
 &= \mathbb{E}[\|(x_i - x_j)P\|_2^2] + \sum_{t=1}^k 4b^2 + 2\mathbb{E}[\langle (x_i - x_j)P, \Delta \rangle] - 2k\sigma^2 \\
 &= \|x_i - x_j\|_2^2 + 4kb^2 + 0 - 2k\sigma^2 \\
 &= \|x_i - x_j\|_2^2 + 4kb^2 + 0 - 2k(2b^2) \\
 &= \|x_i - x_j\|_2^2
 \end{aligned}$$

Note that by the Johnson-Lindenstrauss lemma, we have $\mathbb{E}[\|(x_i - x_j)P\|_2^2] = \|x_i - x_j\|_2^2$. We now show that $2\mathbb{E}[\langle (x_i - x_j)P, \Delta \rangle] = 0$,

Letting $\mathbf{a} = (x_i - x_j)$, we have

$$\begin{aligned}
 2\mathbb{E}[\langle (x_i - x_j)P, \Delta \rangle] &= 2\mathbb{E}[\langle \mathbf{a}P, \Delta \rangle] \\
 &= 2 \sum_{t=1}^k \mathbb{E}[(\mathbf{a}P)_t] \mathbb{E}[\Delta_t] \\
 &= 2 \sum_{t=1}^k \mathbb{E}\left[\left(\sum_{m=1}^d a_m P_{mt}\right)\right] \mathbb{E}[\Delta_t] \\
 &= 2 \sum_{t=1}^k \left(\sum_{m=1}^d a_m \mathbb{E}[P_{mt}]\right) \mathbb{E}[\Delta_t] \\
 &= 0
 \end{aligned}$$

where we have used the independence of P and Δ .

□

Claim 7. Let $S \subset \mathbb{R}^d$ with $|S| = n$. Then, given any two entries in this set $x_i, x_j \in S$, let $y_i = x_i P + \Delta_i$ and $y_j = x_j P + \Delta_j$, where P and Δ_i, Δ_j are the projection matrix and the noise vectors respectively. Then the variance of $\mathcal{D}(y_i, y_j) = \frac{2}{k} \|x_i - x_j\|_2^4 + 2k(7\sigma^4 - \sigma^2) + 4\sigma^2 \|x_i - x_j\|_2^2$, where σ^2 is the variance of the entries of Δ_i and Δ_j .

Proof. Let,

$$\begin{aligned} Z_1 &= \|(x_i - x_j)P\|_2^2 \\ Z_2 &= \|\Delta\|_2^2 \\ Z_3 &= 2\langle (x_i - x_j)P, \Delta \rangle \end{aligned}$$

Then, $\text{Var}(\|x_i P + \Delta_i - x_j P - \Delta_j\|_2^2 - 2k\sigma^2) = \text{Var}(Z_1 + Z_2 + Z_3) - 2k\sigma^2$. Then,

$$\begin{aligned} \text{Var}(Z_1 + Z_2 + Z_3) &= \mathbb{E}[(Z_1 + Z_2 + Z_3)^2] \\ &\quad - (\mathbb{E}[Z_1 + Z_2 + Z_3])^2 \\ &= \mathbb{E}[Z_1^2] - \mathbb{E}[Z_1]^2 + \mathbb{E}[Z_2^2] - \\ &\quad \mathbb{E}[Z_2]^2 + \mathbb{E}[Z_3^2] - \mathbb{E}[Z_3]^2 + \\ &\quad 2\mathbb{E}[Z_1 Z_2] - 2\mathbb{E}[Z_1]\mathbb{E}[Z_2] + \\ &\quad 2\mathbb{E}[Z_2 Z_3] - 2\mathbb{E}[Z_2]\mathbb{E}[Z_3] \\ &\quad + 2\mathbb{E}[Z_1 Z_3] - 2\mathbb{E}[Z_1]\mathbb{E}[Z_3] \\ &= \frac{2}{k} \|x_i - x_j\|_2^4 + 14k\sigma^4 + \\ &\quad 4\sigma^2 \|x_i - x_j\|_2^2 \end{aligned} \tag{1}$$

where we have used the following :

$$\begin{aligned} \mathbb{E}[Z_1^2] - \mathbb{E}[Z_1]^2 &= \frac{2}{k} \|x_i - x_j\|_2^4 \\ \mathbb{E}[Z_2^2] - \mathbb{E}[Z_2]^2 &= 14k\sigma^4 \\ \mathbb{E}[Z_3^2] - \mathbb{E}[Z_3]^2 &= 4\sigma^2 \|x_i - x_j\|_2^2 \end{aligned}$$

Using independence of Z_1 and Z_2 we have $2\mathbb{E}[Z_1 Z_2] = 2\mathbb{E}[Z_1]\mathbb{E}[Z_2]$. For the rest of the variables we have $2\mathbb{E}[Z_2 Z_3] = 2\mathbb{E}[Z_2]\mathbb{E}[Z_3] = 2\mathbb{E}[Z_1 Z_3] = 2\mathbb{E}[Z_1]\mathbb{E}[Z_3] = 0$. Using these, the required result follows. \square

We can consider the probability of the distance recovery algorithm exceeding a fixed error λ . More specifically, letting x_i, x_j be two points in the original space and letting y_i, y_j be the points after the action of the mechanism, we want to know how this value is bounded :

$\Pr [\mathcal{D}(y_i, y_j) - \|x_i - x_j\|_2^2 > \lambda]$. Using the Chebyshev inequality, we get:

$$\begin{aligned} \Pr [|\mathcal{D}(y_i, y_j) - \mathbb{E}[\mathcal{D}(y_i, y_j)]| > \lambda] &\leq \frac{\text{Var}(\mathcal{D}(y_i, y_j))}{\lambda^2} \\ \therefore \Pr [\mathcal{D}(y_i, y_j) - \|x_i - x_j\|_2^2 > \lambda] &\leq \frac{\text{Var}(\mathcal{D}(y_i, y_j))}{\lambda^2} \end{aligned}$$

Therefore, we see that the probability the distance recovery algorithm exceeds a fixed error is proportional to the distance between the original points and the dimensionality of the subspace we project the data into.

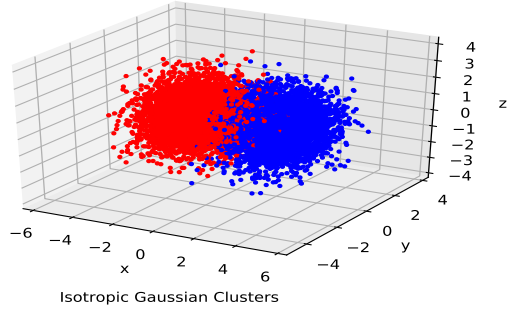


Figure 1. Scatter plot of the original 3D dataset with colors separating the two clusters

5. Experiments

All of our experiments are based on a synthetic dataset comprising of two clusters generated through the procedure developed by Guyon (2003), which is commonly referred to as the Madelon dataset. In this data generation procedure, the data points for each cluster are sampled from an isotropic Gaussian distribution. We fix the Euclidean distance between the cluster centers to 4 and use an ϵ of 4 for all our experiments. We also assume that for two neighboring databases (element-wise or row-wise), the norm of their difference is bounded by 1. That is, for two neighboring databases $X, X' \in \mathbb{R}^{n \times d}$, we have $\|X - X'\|_1 \leq 1$. We use the open source implementation of this data generation procedure provided in Scikit as `sklearn.datasets.make_blob` (Pedregosa et al. (2011)).

Since this dataset is decoupled from any specific problem domain and is a two class clustering problem, it allows us to demonstrate the utility of our mechanism with maximum generality. In order to better understand the generated data in higher dimensions, we first generate data in \mathbb{R}^3 and provide its visualization in Figure 1. Next, we illustrate the effect of the element-wise and row-wise privacy mechanisms defined in Algorithm 1, by visualizing the data using $k = 2$ (i.e. projecting it into \mathbb{R}^2 and making it private) in Figure 2.

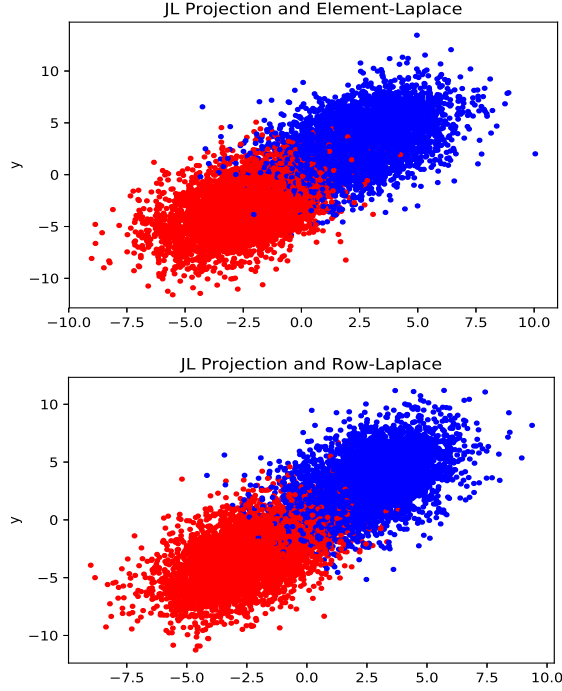


Figure 2. Scatter plots of the 2D dataset for the Element-wise privacy mechanism and the Row-wise privacy mechanism with colors separating the two clusters. Each privacy mechanism turns the spherical data into ellipses that are stretched along the directions of the noise while still maintaining separation.

Using the same dataset defined above, we verify the utility guarantees of our distance recovery algorithm (Algorithm 2). In order to do so, we first sample 1000 pairs of points from the original dataset and calculate the squared Euclidean distance between each pair. This gives us a total of 1000 distances. We then run each pair through our privacy mechanism 1000 times using a new projection and perturbation vector each time (giving us 1 million private pairs). We then use our distance recovery algorithm (Algorithm 2) to calculate the squared Euclidean distance between each private pair of points, giving us 1 million distances for the private data points (1000 distances for each private pair). Next, we plot the distribution of differences in the squared Euclidean distance between the original pair and the private pairs in Figure 3. One can see that the distance recovery algorithm does indeed recover the squared Euclidean distances in expectation. For one run of the experiment, the mean of the differences for element-wise private pairs and row-wise private pairs from the original pairs was found to be 0.006 and -0.011 respectively.

In order to test the utility of our method on the task of clustering, we compare the performance of the k-means clustering algorithm on the original data, element-wise private data, and the row-wise private data. We ran this compari-

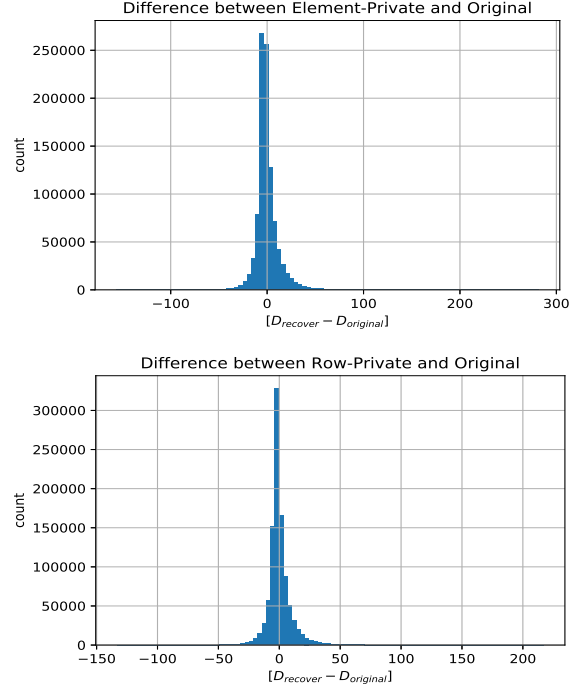


Figure 3. Distribution of the difference between the original squared Euclidean distances and the squared Euclidean distances recovered by our distance recovery algorithm.

son for a number of datasets in which we varied both the original and projected dimensions. The results of this experiments are provided in Table 1. We note that the mechanism provides good utility for smaller k but the utility deteriorates as we increase the dimensionality of the projected subspace, a result that is expected due to the reliance of Laplacian noise on the projection subspace parameter k as shown in Theorems 3, and 5.

We examine the relationship in more detail by plotting the relationship between k and the standard deviation of the data in Figure 4. The formula used for this is $\sqrt{1 + 2b^2}$, where 1 is the variance of the original data and $2b^2$ is the variance of the Laplacian noise. It can be noted that the standard deviation increases with k hence negatively affecting the amount of utility provided by the mechanism. We also note a difference in the standard deviation between element-wise and row-wise privacy mechanisms - row-wise privacy comes at a higher cost utility cost than element-wise privacy.

Our experiments validate the ability of the distance recovery algorithm (Algorithm 2) to recover the squared Euclidean distances between original points and also show that our privacy mechanisms are able to maintain utility in the general task of clustering. We find that the utility of our mechanism deteriorates with an increase in the di-

| Privacy Mechanism | $d=3, k=2$ | $d=10, k=3$ | $d=50, k=10$ | $d=100, k=20$ |
|-------------------|------------|-------------|--------------|---------------|
| None | 0.9783 | 0.9772 | 0.9771 | 0.9797 |
| Element-Wise | 0.9441 | 0.9082 | 0.6954 | 0.6927 |
| Row-Wise | 0.9477 | 0.909 | 0.6796 | 0.6668 |

Table 1. Comparison of performance of k-means clustering between the non-private, element-wise private, and row-wise private data. Here d is the original dimension of the data and k is the projected dimension.

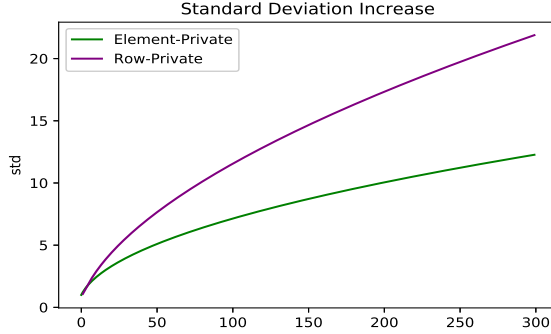


Figure 4. Increase in standard deviation with k

dimensionality of the projection subspace which agrees with Theorems 3, and 5.

6. Related Work

Differential privacy is a framework proposed by (Dwork et al., 2006) that captures precisely how much additional information of an individual is leaked by participating in a database, that would not have been leaked otherwise. There has been extensive research in proposing mechanisms that guarantee differential privacy in the non-interactive setting (Alda & Rubinstein, 2017; Balog et al., 2017; McSherry & Talwar, 2007; Dwork et al., 2014).

(Kenthapadi et al., 2012) developed a randomization mechanism that utilized the JL transform and the Gaussian mechanism (Dwork et al., 2014) to provide non-interactive differential privacy with respect to attribute changes. They showed that their mechanism preserved utility by preserving distances in expectation. However, a shortcoming of this approach was that the privacy guarantees were only provided with respect to attribute changes, and not row level changes, which is a more realistic requirement in practice. Despite that shortcoming, the mechanism was powerful from a privacy perspective, as it had been shown by (Liu et al., 2006) that random projection-based multiplicative perturbation techniques make it impossible to find the exact values of the original data in addition to simply hiding the dimensionality of the data. Further, they showed that if even if the projection matrix is released, the adversary still cannot find the exact value of any elements from

the original data.

(Blocki et al., 2012) showed that the JL transform itself preserved differential privacy and provided utility guarantees in the strict case when only the covariance matrix is released. However, in order to provide privacy guarantees, the data matrix was required to be full rank with eigenvalues above some threshold. Since this is not always feasible in practice, they provided a work around which perturbed all the singular values of the data matrix. In practice, this magnitude of this perturbation can be orders of magnitude larger than the attribute values, hence causing general machine learning algorithms to have extremely poor performance. Along similar lines of using multiplicative random projections to preserve privacy for special problems is the work of (Zhou et al., 2009) who showed that multiplicative random projection methods preserved utility in the case of doing PCA.

Releasing differentially private data raises some fundamental questions about the ability of machine learning algorithms to extract utility from the private data. (Kasiviswanathan et al., 2011) showed that in the PAC learning model with a discrete domain, any finite hypothesis class that is PAC learning is also privately PAC learnable. These results were extended to half space queries by (Blum et al., 2013) and the sample complexities of proper and improper learners were analyzed by (Beimel et al., 2010). However, (Chaudhuri & Hsu, 2011) showed that there exist simple hypothesis classes over continuous domains that have a small VC dimension and for whom it is impossible learn privately with a finite sample size. (Friedman & Schuster, 2010) analyzed the trade-off between privacy, sample complexity, and utility in practice for the case of decision trees.

Another line of research focused on releasing differentially private models with respect to the data (Chaudhuri & Monteleoni, 2009; Evfimievski et al., 2004; Sheffet, 2015a; Zhu et al., 2017). (Chaudhuri et al., 2011) developed a mechanism for private empirical risk minimization that provided private approximates to classifiers and along similar lines (Bassily et al., 2014) analyzed error bounds on such classifiers. Releasing private models also raised questions between the trade-off of privacy and algorithmic complexity, which was analyzed by (Friedman & Schuster, 2010) in practice for the case of de-

cision trees.

7. Conclusions and Future Work

We developed a privacy mechanism that makes it impossible to reconstruct the original data values while also providing utility for general machine learning tasks. We proved privacy guarantees under element and row wise changes, and also proved utility guarantees by proposing an algorithm that maintains pairwise distances between private data points in expectation. We chose the most general task of clustering in order to show the effectiveness of our methods experimentally and validated that it does in fact maintain utility. Noting that the utility of our mechanism deteriorates with an increase in the dimensionality of the projection subspace, we leave open the question of finding a private mechanism that makes reconstruction impossible and provides utility that does not deteriorate with the dimensionality of the problem.

References

- Alda, Francesco and Rubinstein, Benjamin IP. The bernstein mechanism: Function release under differential privacy. In *AAAI*, pp. 1705–1711, 2017.
- Balog, Matej, Tolstikhin, Ilya, and Schölkopf, Bernhard. Differentially private database release via kernel mean embeddings. *arXiv preprint arXiv:1710.01641*, 2017.
- Bassily, Raef, Smith, Adam, and Thakurta, Abhradeep. Differentially private empirical risk minimization: Efficient algorithms and tight error bounds. *arXiv preprint arXiv:1405.7085*, 2014.
- Beimel, Amos, Kasiviswanathan, Shiva Prasad, and Nissim, Kobbi. Bounds on the sample complexity for private learning and private data release. In *Theory of Cryptography Conference*, pp. 437–454. Springer, 2010.
- Bhowmick, Abhishek, Duchi, John, Freudiger, Julien, Kapoor, Gaurav, and Rogers, Ryan. Protection against reconstruction and its applications in private federated learning. *arXiv preprint arXiv:1812.00984*, 2018.
- Blocki, Jeremiah, Blum, Avrim, Datta, Anupam, and Shffet, Or. The johnson-lindenstrauss transform itself preserves differential privacy. In *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*, pp. 410–419. IEEE, 2012.
- Blum, Avrim, Ligett, Katrina, and Roth, Aaron. A learning theory approach to noninteractive database privacy. *Journal of the ACM (JACM)*, 60(2):12, 2013.
- Chaudhuri, Kamalika and Hsu, Daniel. Sample complexity bounds for differentially private learning. In *Proceedings of the 24th Annual Conference on Learning Theory*, pp. 155–186, 2011.
- Chaudhuri, Kamalika and Monteleoni, Claire. Privacy-preserving logistic regression. In *Advances in Neural Information Processing Systems*, pp. 289–296, 2009.
- Chaudhuri, Kamalika, Monteleoni, Claire, and Sarwate, Anand D. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(Mar):1069–1109, 2011.
- Dasgupta, Sanjoy and Gupta, Anupam. An elementary proof of a theorem of johnson and lindenstrauss. *Random Structures & Algorithms*, 22(1):60–65, 2003.
- Dwork, Cynthia and Yekhanin, Sergey. New efficient attacks on statistical disclosure control mechanisms. In *Annual International Cryptology Conference*, pp. 469–480. Springer, 2008.
- Dwork, Cynthia, McSherry, Frank, Nissim, Kobbi, and Smith, Adam. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pp. 265–284. Springer, 2006.
- Dwork, Cynthia, Roth, Aaron, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Dwork, Cynthia, Smith, Adam, Steinke, Thomas, and Ullman, Jonathan. Exposed! a survey of attacks on private data. *Annual Review of Statistics and Its Application*, 4: 61–84, 2017.
- Evmimievski, Alexandre, Srikant, Ramakrishnan, Agrawal, Rakesh, and Gehrke, Johannes. Privacy preserving mining of association rules. *Information Systems*, 29(4):343–364, 2004.
- Friedman, Arik and Schuster, Assaf. Data mining with differential privacy. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 493–502. ACM, 2010.
- Kasiviswanathan, Shiva Prasad, Lee, Homin K, Nissim, Kobbi, Raskhodnikova, Sofya, and Smith, Adam. What can we learn privately? *SIAM Journal on Computing*, 40(3):793–826, 2011.
- Kenthapadi, Krishnaram, Korolova, Aleksandra, Mironov, Ilya, and Mishra, Nina. Privacy via the johnson-lindenstrauss transform. *arXiv preprint arXiv:1204.2606*, 2012.

- Liu, Kun, Kargupta, Hillol, and Ryan, Jessica. Random projection-based multiplicative data perturbation for privacy preserving distributed data mining. *IEEE Transactions on knowledge and Data Engineering*, 18(1):92–106, 2006.
- McSherry, Frank and Talwar, Kunal. Mechanism design via differential privacy. In *Foundations of Computer Science, 2007. FOCS'07. 48th Annual IEEE Symposium on*, pp. 94–103. IEEE, 2007.
- Rahman, Md Atiqur, Rahman, Tanzila, Laganriere, Robert, Mohammed, Noman, and Wang, Yang. Membership inference attack against differentially private deep learning model.
- Sheffet, Or. Differentially private ordinary least squares. *arXiv preprint arXiv:1507.02482*, 2015a.
- Sheffet, Or. Private approximations of the 2nd-moment matrix using existing techniques in linear regression. *arXiv preprint arXiv:1507.00056*, 2015b.
- Shokri, Reza, Stronati, Marco, Song, Congzheng, and Shmatikov, Vitaly. Membership inference attacks against machine learning models. In *Security and Privacy (SP), 2017 IEEE Symposium on*, pp. 3–18. IEEE, 2017.
- Upadhyay, Jalaj. Differentially private linear algebra in the streaming model. *arXiv preprint arXiv:1409.5414*, 2014.
- Zhou, Shuheng, Ligett, Katrina, and Wasserman, Larry. Differential privacy with compression. In *Information Theory, 2009. ISIT 2009. IEEE International Symposium on*, pp. 2718–2722. IEEE, 2009.
- Zhu, Tianqing, Li, Gang, Zhou, Wanlei, and Philip, S Yu. Differentially private deep learning. In *Differential Privacy and Applications*, pp. 67–82. Springer, 2017.