# Differentially Private Mixture of Generative Neural Networks

Gergely Acs<sup>1</sup>, Luca Melis<sup>2</sup>, Claude Castelluccia<sup>3</sup>, and Emiliano De Cristofaro<sup>2</sup> <sup>1</sup>CrySyS Lab, BME-HIT, <sup>2</sup>University College London, <sup>3</sup>INRIA

Abstract-Over the past few years, an increasing number of applications of generative models have emerged that rely on large amounts of contextually rich information about individuals. Owing to possible privacy violations of individuals whose data is used to train these models, however, publishing or sharing generative models is not always viable. In this paper, we introduce a novel solution geared for privately releasing generative models as well as entire high-dimensional datasets produced by these models. We model the generator distribution of the training data by a mixture of k generative neural networks. These networks are trained together, and collectively learn the generator distribution of the given dataset. More specifically, the data is first divided into k clusters using a novel differential private kernel kmeans, then each cluster is given to a separate generative neural network, such as Restricted Boltzmann Machines or Variational Autoencoders, which are trained only on their own cluster using differentially private gradient descent. As the components of our model are neural networks, it can characterize complicated data distributions, and applies to various types of data. We evaluate our approach using the MNIST dataset and a large Call Detail Records (CDR) dataset, and show that it produces realistic synthetic samples, which can also be used to accurately compute arbitrary number of counting queries.

## I. INTRODUCTION

Generative models are an emerging area in machine learning, as recent advances enable the artificial generation of various kinds of data, including images, videos, texts, music. These models are used in a plethora of applications, such as compression [41], denoising [5], inpainting [49], superresolution [26], semi-supervised learning [36], clustering [42], and deep neural networks pretraining [18] in cases where labeled data is expensive. The main goal of generative models is to estimate the underlying distribution of the data and then randomly generate realistic samples according to their estimated distribution. The real distribution-generating data is described with significantly fewer parameters than the number of available samples from this distribution. This "enforced compression" incentivizes the model to describe general features of the training data. Ideally, such generalization should prevent the model to learn any individual-specific information. However, common learning algorithms do not provide such privacy guarantees and often overfit on specific training samples by implicitly memorizing them. For example, model inversion attacks [19] show how an adversary can use a trained model to make predictions of unintended (sensitive) attributes used as input to the model. Therefore, even if only internal model parameters are released, they might still pose significant threats to the privacy of individuals whose data is used for training.

In this paper, we present a novel approach supporting the release of *generative* models while guaranteeing differential privacy to individuals whose data are used to train these models. Differential privacy gained momentum both in the research community and in industry as it provides solid and

measurable privacy guarantees independent of any auxiliary information of the adversary. While previous work explored the use of differential privacy in different areas of machine learning, including deep learning [1], [31], [38], the privacy of generative models has fallen beyond the scope so far.

Generative models play an important role whenever entities holding rich personal datasets are willing or compelled to publish their data, e.g., aiming to monetize it or allow third parties with the appropriate expertise to analyze it. For instance, Call Detail Records (CDRs) collected by telecommunication companies are not only useful to capture interactions between customers, but also to understand their behavior, e.g., for infectious disease spreading or migration patterns.<sup>1</sup> As a result, telcos are often interested in releasing them-more specifically, rather than only releasing specific aggregate statistics, such as certain counting queries or histograms, the ultimate goal is to share an "anonymized" dataset, which replaces the original data in any, perhaps privacy-sensitive, data analytics. Traditional anonymization models, such as k-anonymity, are known to fail on high dimensional data, providing poor utility with insufficient privacy guarantees [3]. A more promising approach is to model the data generating distribution by training a generative model on the original data, and only publish the model along with its (differential private) parameters. Provided with this privacy-preserving model, anybody can generate a synthetic dataset resembling the original (training) data as much as possible without violating differential privacy. The intuition is that generative models have the potential to automatically learn the general features of a dataset including complex regularities such as the subtle and valuable correlation among different attributes.

In this work, we propose a generative model that is a mixture of k generative artificial neural networks (ANNs). These ANNs are trained together and collectively learn the generator distribution of the given dataset. The data is first divided into k clusters using a differential private clustering approach, then each cluster is given to a separate generative neural network, such as Restricted Boltzmann Machines (RBM) [20] or Variational Autoencoders (VAE) [24], which are trained only on their own cluster using differential private gradient descent. A high-level overview of our proposal is shown in Fig. 1. As the components of our model are neural networks, it can characterize complicated data distributions, and potentially applies to various types of data.

Training distinct generative models on different partitions of the dataset has several benefits. First, multiple models can generate more accurate synthetic samples than a single model trained on the whole dataset, as each ANN is trained only on similar data samples. This prevents the mixture model to generate unrealistic synthetic samples which may arise

<sup>&</sup>lt;sup>1</sup>See, e.g., http://www.flowminder.org.



Fig. 1: Overview of our differentially private generative model (DPGM).

from the implausible combination of multiple very different clusters. This scenario is much more likely when the training is perturbed to guarantee differential privacy. Second, each ANN models a different component of the generator distribution, and hence learn any specifics of a cluster faster than a single model. In other words, a single model would need more training epochs than a mixture of generative models to achieve a comparably rich representation of the clusters. As each iteration of the learning algorithm requires some perturbation to guarantee privacy, a mixture model needs less noise which eventually yields more accurate model parameters. Finally, separate models provide larger control over synthetic data generation; unlike a single model, our approach allows to choose the component model and generate synthetic data only from the chosen cluster.

For the data clustering, we use a novel differential private kernel k-means algorithm. Kernel k-means [37] is a nonlinear extension of the classical k-means algorithm and has been shown to be equivalent with most other kernel based clustering algorithms [15]. We first transform the data into a lowdimensional space using random Fourier features [33], and then apply a differential private version of Lloyd's algorithm [7] to find the clusters in the data. Random Fourier features does not only make kernel k-means scalable for large datasets [14], but, unlike standard k-means [7], require to add limited amount of noise to guarantee privacy. Finally, when clusters are created, a generative model is trained on each cluster using differential private stochastic gradient descent (SGD), which is a standard learning technique of many generative ANNs. Previous works added constant amount of noise to the gradient update in each SGD iteration to guarantee differential privacy. Instead, we add noise to each gradient update which is tailored to the data. We prove that our complete scheme provides differential privacy by using the moment accountant method, proposed in [1], which allows to quantify the privacy guarantee of the composition of differential private mechanisms (e.g., noisy k-means iterations followed by noisy SGD iterations) much more accurately than previous work [17].

Contributions. This paper makes the following contributions:

 We propose a novel approach relying on generative neural networks to model the data generating distribution of various kinds of data. It provides differential privacy to each individual in the training data, thus, it can be used to effectively "anonymize" and share large high-dimensional datasets with any potentially adversarial third party.

- 2) As part of our model, we design a novel differential private clustering algorithm based on kernel *k*-means, which efficiently clusters high-dimensional large datasets with strong privacy guarantees.
- 3) We improve the differential private gradient descent algorithm recently proposed by Abadi et al. [1] by using a novel adaptive perturbation technique. We adaptively recompute the magnitude of the noise used to perturb the gradient updates in each SGD iteration, which can lead to significant accuracy improvement of the trained model.
- 4) We evaluate our approach on the MNIST dataset [25] as well as a large Call Detail Records (CDR) dataset, and show that our techniques provide realistic synthetic samples which can also be used to accurately compute arbitrary number of counting queries.

## II. PRELIMINARIES

This section reviews concepts used throughout the rest of the paper. We use the following notation:  $\mathbb{I}$  denotes a universe of items (e.g., set of visited locations, pixels in an image, etc.), where  $|\mathbb{I}| = m$ . A dataset  $D \subseteq 2^{\mathbb{I}}$  is the ensemble of all items of some set of individuals. A record, which is a non-empty subset of  $\mathbb{I}$ , refers to all items of an individual from D and is represented by a binary vector  $\mathbf{x}$  of size m.

## A. Restricted Boltzmann Machines (RBM)

A Restricted Boltzmann Machine (RBM) is a bipartite undirected graphical model composed of m visible and ninvisible (or latent) binary random variables denoted by, respectively,  $\mathbf{v} = (v_1, v_2, \dots, v_m)$  and  $\mathbf{h} = (h_1, h_2, \dots, h_n)$ . In our case, visible variables represent the attributes of Dand their values are composed of records from D. Hidden variables capture the dependencies between different visible variables (i.e., dependencies between the items in I). As the above model is a Markov random field with strictly positive joint probability distribution p over the model variables, p can be represented as a Boltzmann distribution defined as:

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}$$
(1)

where  $Z = \sum_{\mathbf{v},\mathbf{h}} e^{-E(\mathbf{v},\mathbf{h})}$  is the partition function,  $E(\mathbf{v},\mathbf{h})$ the energy function, i.e.,  $E(\mathbf{v},\mathbf{h}) = -\sum_{i=1}^{n} \sum_{j=1}^{m} v_{ij}h_iv_j - \sum_{j=1}^{m} b_jv_j - \sum_{i=1}^{n} c_ih_i$ , with  $w_{ij}$  being real valued weights describing the inter-dependency between  $v_j$  and  $h_i$ , and  $b_j, c_i$ real valued bias terms associated with the *j*th visible and *i*th hidden units, respectively. Using matrix notation,  $E(\mathbf{v},\mathbf{h}) = -\mathbf{v}^{\top}\mathbf{W}\mathbf{h} - \mathbf{b}^{\top}\mathbf{v} - \mathbf{c}^{\top}\mathbf{h}$ , where  $\mathbf{W} = [\![w]\!]_{i,j}$ ,  $\mathbf{c} = [c]_i$ , and  $\mathbf{b} = [b]_j$ . The goal is to approximate the true data generating distribution with the Boltzmann distribution *p*, given in Eq. (1). To this end, we train the RBM model on dataset *D* to compute parameters  $\mathbf{W}, \mathbf{c}, \mathbf{b}$ .

There are a few algorithms to train RBMs, that approximate or relate to gradient descent on the log-likelihood of the data. If  $\theta = (\mathbf{W}, \mathbf{b}, \mathbf{c})$ , then we want to maximize the likelihood function  $\mathcal{L}(\theta|D) = \prod_{\mathbf{x}\in D} p(\mathbf{x}|\theta)$  given dataset D, where  $\mathbf{x} \in \{0,1\}^m$  is a record from D and p is the Boltzmann distribution defined in Eq. (1). A numerical approximation, gradient descent, is used where the model parameters  $\theta$  are iteratively updated using D and the gradient of the loglikelihood function as:  $\theta_{t+1} = \theta_t + \eta \frac{\partial \log \mathcal{L}(\theta_t|D)}{\partial \theta_t}$ , with  $\eta \in \mathbb{R}^+$ being the learning rate. The model parameters are updated until the log-likelihood converges. In this paper, we employ Persistent Contrastive Divergence [43].

#### B. Variational Autoencoder (VAE)

A variational autoencoder [24] consists of two neural networks (an encoder and a decoder), and a loss function. The encoder compresses data into a latent space (z) while the decoder reconstructs the data given the hidden representation. Let  $\mathbf{x}$  be a random vector of m observed variables, which are either discrete or continuous. Let z be a random vector of n latent continuous variables. The probability distribution between **x** and **z** assumes the form  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x} \mid \mathbf{z})$ , where  $\theta$  indicates that p is parametrized by  $\theta$ . Also, let  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$ be a recognition model whose goal is to approximate the true and intractable posterior distribution  $p_{\theta}(\mathbf{z} \mid \mathbf{x})$ . We can then define a lower-bound on the log-likelihood of  $\mathbf{x}$  as follows:  $\mathcal{L}(\mathbf{x}) = -D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \mid\mid p_{\theta}(\mathbf{z})) + \mathbf{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})}[\log p_{\theta}(\mathbf{x} \mid \mathbf{z})].$ The first term pushes  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$  to be similar to  $p_{\theta}(\mathbf{z})$  ensuring that, while training, VAE learns a decoder that, at generation time, will be able to invert samples from the prior distribution such they look just like the training data. The second term can be seen as a form of reconstruction cost, and needs to be approximated by sampling from  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$ .

In VAEs, we propagate the gradient signal through the sampling process and through  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$  using the *reparametrization trick*. This is done by making  $\mathbf{z}$  be a deterministic function of  $\phi$  and some noise  $\epsilon$ , i.e.,  $\mathbf{z} = f(\phi, \epsilon)$ . For instance, sampling from a normal distribution can be done like  $\mathbf{z} = \mu + \sigma \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ . The reparametrization trick can be viewed as an efficient way of *adapting*  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$  to help improve the reconstruction. We train the variational autoencoder using stochastic gradient descent to optimize the loss with respect to the parameters of the encoder and decoder  $\theta$  and  $\phi$ .

## C. Kernel k-means with Random Features

Given a set of samples  $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , k-means linearly separates D into k clusters  $C_1, C_2, \dots, C_k$   $(k \le N)$ so that it aims to minimize the error  $\sum_{i=1}^k \sum_{\mathbf{x} \in C_i} ||\mathbf{x} - \mathbf{c}_i||_2^2$ , where  $\mathbf{c}_i = \sum_{\mathbf{x} \in C_i} \mathbf{x} / |C_i|$  is the centroid of cluster  $C_i$ . Although this problem is NP-hard, there are efficient heuristic algorithms (such as Lloyd's algorithm) which iteratively refines clustering and converge quickly to a local optimum. However, k-means can provide very inaccurate clustering of linearly non-separable data, which are very common in practice. To overcome this shortcoming, kernel k-means [37] first maps samples from input space to a higher dimensional feature space through a non-linear transformation  $\Phi$ , then applies standard k-means on  $\{\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_N)\}$ . Hence, kernel kmeans provides linear separators of clusters in feature space which correspond to non-linear separators in input space. Kernel k-means iteratively computes  $||\Phi(\mathbf{x}) - \mathbf{c}'_i||_2^2$  for each sample **x** to decide which cluster a sample belongs to, where  $\mathbf{c}'_i = \sum_{\mathbf{x}\in C_i} \Phi(\mathbf{x})/|C_i|$ . To do so, the inner product  $\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle$  must be known for all  $\mathbf{x}, \mathbf{y} \in D$ . Since  $\Phi(\cdot)$ is hard to explicitly compute due to its large, often infinite dimension, the kernel trick is applied;  $\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle = \kappa(\mathbf{x}, \mathbf{y}),$ where  $\kappa$  is an easily computable kernel function. Still, this approach requires evaluating  $\kappa$  for all pairs of samples and store the results, which is not scalable for large datasets.

To make kernel k-means scalable, the kernel function can be approximated with low-dimensional explicit feature maps. In particular, the samples are first mapped to a low-dimensional Euclidean inner product space using an explicit random feature map  $z : \mathbb{R}^m \to \mathbb{R}^d$  so that  $\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle \approx \langle z(\mathbf{x}), z(\mathbf{y}) \rangle$ . Then, standard k-means is applied on the low-dimensional mapped samples  $\{z(\mathbf{x}_1), z(\mathbf{x}_2), \ldots, z(\mathbf{x}_N)\}$  in  $\mathbb{R}^d$  to approximate the result of the kernel k-means with implicit feature map  $\Phi$  and kernel  $\kappa$ . Although the approximation error decreases as d increases, quite accurate approximations can be obtained even for d < m. Explicit nonlinear feature maps have already been proposed for shift-invariant kernels (e.g., generalized RBF kernels) [44] as well as polynomial kernels [32] among others.

#### D. Differential Privacy (DP)

Differential privacy allows a party to privately release a dataset: using perturbation mechanisms, a function of an input dataset is modified, so that any information which can discriminate a record from the rest of the dataset is bounded [16].

Definition 1 (Privacy loss): Let  $\mathcal{A}$  be a privacy mechanism which assigns a value  $Range(\mathcal{A})$  to a dataset D. The privacy loss of  $\mathcal{A}$  with datasets D and D' at output  $O \in Range(\mathcal{A})$  is a random variable  $\mathcal{P}(\mathcal{A}, D, D', O) = \log \frac{\Pr[\mathcal{A}(D)=O]}{\Pr[\mathcal{A}(D')=O]}$  where the probability is taken on the randomness of  $\mathcal{A}$ .

Definition 2 (( $\epsilon, \delta$ )-Differential Privacy [16]): A privacy mechanism  $\mathcal{A}$  guarantees ( $\varepsilon, \delta$ )-differential privacy if for any database D and D', differing on at most one record, and for any possible output  $S \subseteq Range(\mathcal{A}), Pr[\mathcal{A}(D) \in$  $S] \leq e^{\varepsilon} \times Pr[\mathcal{A}(D') \in S] + \delta$  or, equivalently,  $\Pr_{O \sim \mathcal{A}(D)}[\mathcal{P}(\mathcal{A}, D, D', O) > \varepsilon] \leq \delta$ .

This definition guarantees that every output of algorithm  $\mathcal{A}$  is almost equally likely (up to  $\varepsilon$ ) on datasets differing in a single record except with probability at most  $\delta$ , preferably smaller than 1/|D|. Intuitively, this guarantees that an adversary, provided with the output of  $\mathcal{A}$ , can draw almost the same conclusions about any individual no matter if this individual is included in the input of  $\mathcal{A}$  or not [16].

Differential privacy maintains composition, i.e., if each of  $\mathcal{A}_1, \ldots, \mathcal{A}_k$  is  $(\varepsilon, \delta)$ -DP, then their k-fold adaptive composition<sup>2</sup> is  $(k\varepsilon, k\delta)$ -DP. However, a tighter upper bound can be derived on the privacy loss of the composite using a generic Chernoff bound. In particular, it follows from Markov's inequality that  $\Pr[\mathcal{P}(\mathcal{A}, D, D', O) \ge \varepsilon] \le$  $\mathbb{E}[\exp(\lambda \mathcal{P}(\mathcal{A}, D, D', O))]/\exp(\lambda\varepsilon)$  for any output  $O \in$  $Range(\mathcal{A})$  and  $\lambda > 0$ . This implies that  $\mathcal{A}$  is  $(\varepsilon, \delta)$ -DP (see Definition 2) with  $\delta = \min_{\lambda} \exp(\alpha_{\mathcal{A}}(\lambda) - \lambda\varepsilon)$ , where  $\alpha_{\mathcal{A}}(\lambda) = \max_{D,D'} \log \mathbb{E}_{O \sim \mathcal{A}(D)}[\exp(\lambda \mathcal{P}(\mathcal{A}, D, D', O))]$  is the log of the moment generating function of the privacy loss.

Theorem 1 (Moments accountant [1]): Let  $\alpha_{\mathcal{A}_i}(\lambda)$  be  $\max_{D,D'} \log \mathbb{E}_{O \sim \mathcal{A}(D)}[\exp(\lambda \mathcal{P}(\mathcal{A}, D, D', O))]$  and  $\mathcal{A}_{1:k}$  the *k*-fold adaptive composition of  $\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_k$ . It holds:

1) 
$$\alpha_{\mathcal{A}_{1:k}}(\lambda) \leq \sum_{i=1}^{k} \alpha_{\mathcal{A}_i}(\lambda)$$

2)  $\mathcal{A}_{1:k}$  is  $(\varepsilon, \min_{\lambda} \exp(\sum_{i=1}^{k} \alpha_{\mathcal{A}_i}(\lambda) - \lambda \varepsilon))$ -differential private

where  $A_1, A_2, \ldots, A_k$  use independent coin tosses.

There are a few ways to achieve DP, including the Gaussian mechanism [16]. A fundamental concept of all of them is the *global sensitivity* of a function [16].

Definition 3 (Global  $L_p$ -sensitivity): For any function  $f : \mathcal{D} \to \mathbb{R}^d$ , the  $L_p$ -sensitivity of f is  $\Delta_p f = \max_{D,D'} ||f(D) - f(D')||_p$ , for all D, D' differing in at most one record, where  $|| \cdot ||_p$  denotes the  $L_p$ -norm.

**Gaussian Mechanism.** The Gaussian Mechanism (GM) [16] consists of adding Gaussian noise to the true output of a function. In particular, for any function  $f : \mathcal{D} \to \mathbb{R}^d$ , GM is defined as  $\mathcal{G}(D) = f(D) + [\mathcal{N}_1(0, \Delta_2 f \cdot \sigma), \dots, \mathcal{N}_d(0, \Delta_2 f \cdot \sigma)]$ , where  $\mathcal{N}_i(0, \Delta_2 f \cdot \sigma)$  are i.i.d. normal random variables with zero mean and variance  $(\Delta_2 f \cdot \sigma)^2$ .

Lemma 1: 
$$\alpha_{\mathcal{G}}(\lambda) = (\lambda^2 + \lambda)/4\sigma^2$$
  
Proof: See Appendix A.

Given  $\alpha_{\mathcal{G}}(\lambda)$ , the exact privacy cost  $\varepsilon$  (or  $\delta$ ) of the k-fold adaptive composition of  $\mathcal{G}$  is computed based on Theorem 1.

## III. DPGM: DIFFERENTIALLY PRIVATE GENERATIVE MODEL

In this section, we present our Differential Private Generative Model (DPGM), which is described in Alg. 1 and also illustrated in Fig. 1. The dataset D is first partitioned into k clusters, denoted by  $\hat{D}_1, \hat{D}_2, \ldots, \hat{D}_k$ , which are in turn used to train k distinct generative models, where the parameters of the resulting models are denoted, respectively, by  $\theta_1, \theta_2, \ldots, \theta_k$ . Data samples are similar within a cluster, thus, generative models simultaneously trained on each partition converge faster than a single model trained on the whole dataset D. As  $\theta_1, \theta_2, \ldots, \theta_k$  are learnt using perturbed gradient descent, they can be released and used to generate synthetic data using the k generative models.

Our learning approach involves two main steps: (1) records in D are clustered in a random feature space using differential private kernel k-means (see Section III-A) into clusters

## Algorithm 1: DPGM: Differentially Private Generative Model

**Input:** Dataset:  $D = {\mathbf{x}_1, ..., \mathbf{x}_N}$ , # of custers: k, k-means iterations:  $T_{\mathcal{K}}$ , SGD iterations:  $T_{\mathcal{S}}$ , Noise scales:  $\sigma_{\mathcal{C}}, \sigma_{\mathcal{K}}, \sigma_{\mathcal{S}}$ **1** Cluster data records in D:

$$\{\hat{D}_1, \hat{D}_2, \dots, \hat{D}_k\} = \text{DPkmeans}(k, T_{\mathcal{K}}, D, \sigma_{\mathcal{C}}, \sigma_{\mathcal{K}})$$

2 Initialize 
$$\theta_1, \theta_2, \ldots, \theta_k$$
 randomly

3 for  $t \in [T_S]$  do

4 Select  $(\hat{D}_s, \theta_s) \in \{(\hat{D}_1, \theta_1), \dots, (\hat{D}_k, \theta_k)\}$  with probability  $|\hat{D}_s|/|D|$ 

5 Update parameters of model  $\theta_s$ :

6  $\theta_s = \text{DP-SGD}(\hat{D}_s, \theta_s, \sigma_c, \sigma_s) //\text{see Alg. 4}$ Output:  $\theta_1, \theta_2, \dots, \theta_k$ 

Algorithm 2: DPkmeans: Private kernel k-means with Random Fourier
Features
<b>Input:</b> Data: $D = {\mathbf{x}_1, \dots, \mathbf{x}_N}$ , Cluster number: k, Iterations: T,
Feature number: d, Kernel function: $\kappa$ , Noise scales: $\sigma_{\mathcal{C}}, \sigma_{\mathcal{K}}$
<b>1</b> Compute Features: $\mathbf{w}_i \sim_{\text{iid}} p(\mathbf{w})$ for all $1 \leq i \leq d$ , where

 $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{w}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{w}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{w}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{w}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{w}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \kappa(\mathbf{w}) d\mathbf{w}$   $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{w} \rangle) \pi(j\langle \mathbf{w}, \mathbf{w} \rangle) \pi(j$ 

 $\hat{D}_1, \hat{D}_2, \ldots, \hat{D}_k$ ; (2) a generative model (e.g., RBM [20] or VAE [24]) with parameter  $\theta_i$  is trained on cluster  $\hat{D}_i$  (see Section III-B) using differential private gradient descent, where the training data are composed of the records of  $\hat{D}_i$ . In each SGD iteration (Line 4-6 in Alg. 1), a model  $\theta_s$  is chosen uniformly at random along with its corresponding training data  $\hat{D}_s$ , and a single SGD iteration is performed to update  $\theta_s$  using a random sample S of  $\hat{D}_s$  with size L (Line 7 in Alg. 1). The output of our algorithm are composed of the parameters of the trained generative models, i.e.,  $\theta_1, \theta_2, \ldots, \theta_k$ . Finally, these privately trained k models can be used to generate synthetic records which resemble the original ones, i.e., preserve their general characteristics that are not specific to any single individual (up to  $\varepsilon$  and  $\delta$  which is computed in Section III-D).

## A. Private kernel k-means

We now discuss our private kernel k-means algorithm, presented in Alg. 2. It first transforms the data D into a lowdimensional representation  $D' = \{z(\mathbf{x}_1), \ldots, z(\mathbf{x}_N)\}$  using randomized Fourier feature map  $z : \mathbb{R}^m \to \mathbb{R}^d$  [33], and then applies standard differential private k-means [7] on these lowdimensional features.

Specifically, 
$$z : \mathbb{R}^m \to \mathbb{R}^d$$
 is defined as:  
 $z(\mathbf{x}) = \sqrt{\frac{2}{d}} \left[ \cos(\langle \mathbf{w}_1, \mathbf{x} \rangle + b_1), \dots, \cos(\langle \mathbf{w}_d, \mathbf{x} \rangle + b_d) \right]$  (2)

where each  $\mathbf{w}_i \in \mathbb{R}^m$  is drawn independently from  $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j\langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$ , i.e.,  $p(\mathbf{w})$  is the Fourier trans-

<sup>&</sup>lt;sup>2</sup>The output of  $A_{i-1}$  is used as an input of  $A_i$ , i.e., their executions are not necessarily independent except their coin tosses.

form of kernel function  $\kappa$ , and  $b_i \in \mathbb{R}$  is chosen from  $[0, 2\pi)$ uniformly at random. In particular, Bochner's theorem implies that  $p(\mathbf{w})$  is a valid probability density function, if  $\kappa$  is continuous, positive-definite, and shift-invariant kernel. Hence,  $\kappa(\mathbf{x}, \mathbf{y}) = \kappa(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{R}^m} \exp(j\langle \mathbf{w}, \mathbf{x} - \mathbf{y} \rangle) p(\mathbf{w}) d\mathbf{w} =$  $\mathbb{E}_{\mathbf{w}, b}[\langle \sqrt{2}\cos(\langle \mathbf{w}, \mathbf{x} \rangle + b), \sqrt{2}\cos(\langle \mathbf{w}, \mathbf{y} \rangle + b) \rangle]$ , where the expectation is approximated with the empirical mean over drandomly chosen values of  $\mathbf{w}$  and b [33].

Standard DP k-means [7] releases the noisy cluster centers which are computed iteratively using a noisy variant of Lloyd's algorithm; in each iteration, gaussian noise with scale  $\sqrt{2\sigma_{\mathcal{K}}}$ is added to the size of all clusters, and with scale  $\sqrt{2}\sigma_{\mathcal{K}}C_s$ to the sum of all cluster members in each cluster. These noisy values are used to compute the noisy cluster centers  $\{\hat{\mathbf{c}}_1,\ldots,\hat{\mathbf{c}}_k\}^3$ . To determine the scale of the gaussian noise, the  $L_2$ -sensitivity of the cluster size and that of the sum of norms must be known within each cluster. Although the  $L_2$ sensitivity of the set of cluster size is always  $\sqrt{2}$  (a single record can change the size of at most 2 clusters), such a priori bound does not exist for the  $L_2$ -norm of the feature vectors in general. Hence, we need to clip all feature vectors in  $L_2$ -norm before applying standard DP k-means, where the clipping threshold  $C_s$  should be set to the average norm of the feature vectors (i.e.,  $(1/N) \sum_{\mathbf{x} \in D} ||z(\mathbf{x}_i)||_2$ ) and is approximated by Alg. 3 (see Section III-C). Replacing  $z(\mathbf{x}_i)$ with  $\hat{z}(\mathbf{x}_i) = z(\mathbf{x}_i) / \max(1, ||z(\mathbf{x}_i)||_2 / C_s)$  guarantees that all feature vectors are kept as long as their norm is less then  $C_s$ , otherwise they are scaled down to have a norm of  $C_s$ .

Nevertheless, for some kernel functions, such as the Radial Basis Function (RBF), a small norm bound  $C_s$  can be computed analytically – see Theorem 2. Interestingly, this bound is constant for any input data and feature size independently of the width  $\gamma$  of the RBF kernel. Therefore, as opposed to standard k-means [7], our approach can detect linearly non-separable clusters, and, used with RBF kernel, add constant noise to feature vectors independently of their size d.

Theorem 2: If  $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-\gamma ||\mathbf{x} - \mathbf{y}||_2)$ , then  $\mathbb{E}[||z(\mathbf{x})||_2] \leq 1$  for any  $\mathbf{x} \in \{0, 1\}^*$  and  $\gamma$ , where the expectation is taken on the randomness of z.

Therefore, DP kernel k-means has two main advantages over standard DP k-means [7]. First, kernel k-means can find linearly non-separable clusters. Second, if it is used with RBF kernel, the added noise is independent of the  $L_2$ -norm of the data records. As we show in Section IV, this can lead to much larger clustering accuracy especially for stringent privacy requirements (i.e., for  $\varepsilon < 0.5$ ) even for large dimensional data.

## B. Private Stochastic Gradient Descent

We now present our private SGD technique—Alg. 4 outlines a single SGD iteration. Our starting point is the work by Abadi et al. [1]: similar to theirs, our solution provides differential privacy to the training data by first clipping the

#### Algorithm 3: DPNorm: Private Approximation of Average Norm

Input: Data:  $S = \{\mathbf{x}_{c_1}, \dots, \mathbf{x}_{c_{\lfloor S \rfloor}}\}$ , Noise scale:  $\sigma_{\mathcal{C}}$ , Max. norm bound:  $C_{\max}$ , Max. number of discretized norm bounds: w1  $C_j \leftarrow j \cdot C_{\max}/w$  for  $0 \le j \le w$ 

2  $C_s \leftarrow \arg\max_{j\geq 1}\{t_j + \mathcal{N}(0, \sqrt{2}\sigma_{\mathcal{C}})\}, \text{ where } t_j = |\{\mathbf{x} \in S : C_{j-1} < ||\mathbf{g}(\mathbf{x})||_2 \le C_j\}|$ Output:  $C_s$ 

Algorithm	4:	Private	Stochastic	Gradient	Descent

 $\begin{array}{l} \text{Input: } Data: \ \hat{D}, \ Model \ parameters: \ weights \ and \ biases \ \theta, \ Noise \\ scales: \ \sigma_{\mathcal{C}}, \ \sigma_{\mathcal{S}}, \ Loss \ function: \ \mathcal{L}(\theta) = \frac{1}{|\hat{D}|} \sum_{i} \mathcal{L}(\theta, \mathbf{x}_{c_{i}}), \\ \text{Learning rate: } \eta, \ \text{Batch size: } L \\ \textbf{1 } \textbf{Sampling: } Take \ a \ random \ sample \ S = \{\mathbf{x}_{c_{1}}, \ldots, \mathbf{x}_{c_{|\mathcal{S}|}}\} \ \text{of } \hat{D} \ \text{with} \\ \text{sampling probability } q = L/|\hat{D}| \\ \textbf{2 } \textbf{Compute Gradient: For each } \mathbf{x}_{c_{i}} \in S, \ \text{compute} \\ \mathbf{g}(\mathbf{x}_{c_{i}}) \leftarrow \nabla_{\theta} \mathcal{L}(\theta, \mathbf{x}_{c_{i}}) \\ \textbf{3 } \textbf{Clip Gradient: } S' \leftarrow \{\mathbf{g}(\mathbf{x}_{c_{i}}), \ldots, \mathbf{g}(\mathbf{x}_{c_{|\mathcal{S}|}})\} \\ \textbf{4 } \ C_{s} \leftarrow \text{DPNorm}(S', \sigma_{\mathcal{C}}) \ // \text{see Alg. } 3 \\ \textbf{5 } ; \ \hat{\mathbf{g}}(\mathbf{x}_{c_{i}}) \leftarrow \mathbf{g}(\mathbf{x}_{c_{i}})/\max\left(1, \frac{||\mathbf{g}(\mathbf{x}_{c_{i}})||_{2}}{C_{s}}\right) \\ \textbf{6 } \ \textbf{Add noise: } \ \tilde{\mathbf{g}} \leftarrow \frac{1}{L} \left(\sum_{i=1}^{|\mathcal{S}|} \hat{\mathbf{g}}(\mathbf{x}_{c_{i}}) + \mathcal{N}(0, \sqrt{2}\sigma_{\mathcal{S}}C_{s}\mathbf{I})\right) \\ \textbf{7 } \ \textbf{Descent: } \theta \leftarrow \theta - \eta \ \tilde{\mathbf{g}} \\ \mathbf{Output: } \theta \end{array}$ 

norm of the gradient update of each record, and then perturbing these clipped gradients by the Gaussian mechanism. However, we achieve better accuracy as the clipping threshold is selected adaptively in each SGD iteration. In particular, in each SGD iteration, we also (1) compute the gradient of the loss function  $\mathcal{L}$  on a random subset S of records (denoted as "batch") in Line 2 of Alg. 4, (2) clip the  $L_2$  norm of the gradient of each record in S to have a norm at most  $C_s$  (in Line 3), (3) add gaussian noise  $\mathcal{N}(0, \sqrt{2\sigma_s C_s I})$  to the average of these clipped gradient updates (Line 6), and finally (4) perform the descent step (Line 7). At the end, the updated model parameters  $\theta$  are returned. A complete training epoch on the whole dataset D consists of (|D|/L) SGD iterations, which are required to process all records in every cluster on average. Indeed, each record in a cluster  $\hat{D}_s$  is selected with probability  $(|\hat{D}_s|/\sum_{i=1}^k |\hat{D}_i|) \times (L/||\hat{D}_s|) = L/|D|$ , where  $\sum_{i=1}^k |\hat{D}_i| = |D|$ . Notice that the  $L_2$ -sensitivity of  $\sum_i \hat{\mathbf{g}}(\mathbf{x}_i)$ is  $\sqrt{2C_s}$ , as the norm of every  $\hat{\mathbf{g}}(\mathbf{x}_i)$  is at most  $C_s$ , and one record can change at most two clusters.

## C. Adaptive selection of the norm bound

Both our private SGD method (in Line 4 of Alg. 4) and private kernel k-means (in Line 4 of Alg. 2) require the differential private computation of the average  $L_2$ -norm in a given set of records, which is then used as the clipping threshold  $C_s$  in both algorithms. For this purpose, these algorithms invoke DPNorm which is detailed in Alg. 3. In fact, our SGD technique differs from the original private SGD method [1] in the selection of the norm bound  $C_s$  (in Line 3-5 of Alg. 4). In the original approach [1],  $C_s$  is provided as input to the private SGD and no guideline is given how to compute its value without violating differential privacy. Moreover, the selection of the norm bound  $C_s$  has a large impact on the performance of the private SGD in general. If  $C_s$ is too small, there will be slow convergence. Conversely, if it is too large, unnecessarily large gaussian noise will be introduced

<sup>&</sup>lt;sup>3</sup>At the beginning, we initialize the clusters centers to random records drawn from *publicly available* non-sensitive data which are generated by the same distribution as the sensitive data. We only need k representative samples, and such public datasets already exist for images, location data, or even medical data.

on the gradient update. Intuitively,  $C_s$  should be adjusted so that  $||\mathbf{g}(\mathbf{x}_{c_i})||_2 \approx C_s$  for each record  $\mathbf{x}_{c_i}$ . This guarantees that the contribution of  $\mathbf{x}_{c_i}$  to  $\tilde{\mathbf{g}}$  is maximally preserved with the smallest relative error. Hence, instead of fixing  $C_s$ for the whole training, we aim to compute  $C_s$  adaptively for each batch as  $C_s = (1/L) \sum_i ||\mathbf{g}(\mathbf{x}_{c_i})||_2$ . This adaptive approach would ensure fast convergence with small error, and also adapt to the gradient update of every batch. Indeed, SGD is iterative, so the gradient update  $\tilde{\mathbf{g}}$  of a batch/iteration depends on that of the previous batch/iteration, which means that  $(1/L) \sum_i ||\mathbf{g}(\mathbf{x}_{c_i})||_2$  is different for each batch.

In DPNorm (see Alg. 2), the computation of the average norm in a set S of records is randomized to guarantee privacy. A naive solution is to add Gaussian noise to this average, i.e.,  $C_s = (1/|S|) \sum_{\mathbf{x} \in S} ||\mathbf{x}||_2 + \mathcal{N}(0, s \cdot \sigma'/L),$ where  $s \geq \max_{\mathbf{x} \in S} ||\mathbf{x}||_2$ . However,  $\max_{\mathbf{x} \in S} ||\mathbf{x}||_2$  is datadependent and can be too large if there are outliers in S. Instead, we approximate  $C_s$  such that its value is close to the norm of many records in S, i.e., it is a good approximator of  $(1/L)\sum_{\mathbf{x}\in S}||\mathbf{x}||_2$ . In particular, we discretize the domain of  $C_s$  by dividing  $(0, C_{\text{max}})$  uniformly into w intervals (Line 1) of Alg. 3). Then, we use the Gaussian mechanism (Line 2 of Alg. 3) to select among the upper bounds  $C_j = jC_{\max}/w$ of these intervals  $(0 \le j \le w)$ , which will be the norm bound  $C_s$  for S. Specifically, we build a histogram where bin *i* equals the number of records whose gradient norm falls within  $(C_{i-1}, C_i]$ . Then, the (noisy) mode of this histogram is computed by adding independent gaussian noise  $\mathcal{N}(0, \sqrt{2\sigma_c})$ to each count, and selecting the bin which has the greatest noisy count. Note that the  $L_2$ -sensitivity of the histogram is always bounded by  $\sqrt{2}$  no matter how large  $\max_{\mathbf{x}\in S} ||\mathbf{x}||_2$  is.

## D. Privacy Analysis

DPGM is the composition of private kernel k-means and private SGD. Let  $\mathcal{K}$  denote the private kernel k-means algorithm whose output is the noisy mapped cluster centers after  $T_{\mathcal{K}}$  clustering iterations (i.e.,  $\mathcal{K}(D) = \{\hat{\mathbf{c}}_1, \dots, \hat{\mathbf{c}}_k\}$ ).  $\mathcal{K}$  is composed of (1) selecting the norm bound using DPNorm and (2)  $T_{\mathcal{K}}$  iterations of k-means. Let  $\mathcal{G}_1$  denote the gaussian mechanism which selects the norm bound as per Section III-C. A single k-means iteration is the 2-fold adaptive composition of two gaussian mechanisms  $\mathcal{G}_2$  and  $\mathcal{G}_3$  (in Line 10-11 of Alg. 2), where  $\mathcal{G}_2$  perturbs the cluster size (Line 10), while  $\mathcal{G}_3$  adds noise to the sum of Fourier features of the cluster members (Line 11). The  $L_2$ -sensitivity of the size of every clusters is  $\sqrt{2}$ , as changing a single record can change the size of at most two clusters. Similarly, the  $L_2$ -sensitivity of the sum of Fourier features of the cluster members is  $\sqrt{2C_s}$ as it is detailed in Section III-A. Since  $\mathcal{K}$  is the  $T_{\mathcal{K}}$ -fold adaptive composition of  $T_{\mathcal{K}}$  clustering iterations, it follows from Theorem 1 and Lemma 1 that

$$\alpha_{\mathcal{K}}(\lambda) \le T_{\mathcal{K}}(\alpha_{\mathcal{G}_1}(\lambda) + \alpha_{\mathcal{G}_2}(\lambda) + \alpha_{\mathcal{G}_3}(\lambda))$$
  
$$\le T_{\mathcal{K}}(\lambda^2 + \lambda)(1/4\sigma_{\mathcal{C}}^2 + 1/2\sigma_{\mathcal{K}}^2)$$
(3)

Note that if the RBF kernel is used in kernel k-means (i.e.,  $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-\gamma ||\mathbf{x} - \mathbf{y}||_2)$  in Alg. 2), then  $\alpha_{\mathcal{G}_1}(\lambda) = 0$  and  $\alpha_{\mathcal{K}}(\lambda) \leq T_{\mathcal{K}}(\lambda^2 + \lambda)/2\sigma_{\mathcal{K}}^2$  since  $C_s = 1$  is a priori bound on the  $L_2$ -norm of every feature vector (cf. Theorem 2).

Let  $S_k$  denote the private SGD algorithm whose output is the noisy model parameters after  $T_S$  SGD iterations (i.e.,  $S(D) = \{\theta_1, \ldots, \theta_k\})^4$ , and the input is the cluster centers  $\{\hat{c}_1, \ldots, \hat{c}_k\}$  provided by  $\mathcal{K}$ . At the very beginning, S assigns each record to its closest cluster center in feature space to obtain k non-overlapping training sets (this is implemented by the last iteration in Alg. 2). Changing a single record alters at most a single record in at most 2 training sets (clusters), as the modified record can be moved from one to another training set. Since all training sets are non-overlapping, each record is selected in an SGD iteration with probability  $q = (|\hat{D}_s|/|D|) \times (L/||\hat{D}_s|) = L/|D|$  for any k. Moreover, each of the k models are trained independently, so  $\alpha_{\mathcal{S}_k}(\lambda) \leq \alpha_{\mathcal{S}_1}(\lambda)$ , where  $\mathcal{S}_1$  denotes the case when k = 1 (i.e., a single model is trained on the whole dataset D during  $T_S$  epochs).

The complete SGD training of  $S_1$  (Line 3-6 in Alg. 1) is the  $T_{\mathcal{S}}$ -fold adaptive composition of  $T_{\mathcal{S}}$  SGD iterations, where we jointly use two perturbation mechanisms  $\mathcal{G}_4$  and  $\mathcal{G}_5$  in each iteration;  $\mathcal{G}_4$  selects a batch uniformly at random and computes the norm bound  $C_s$  (Line 4 of Alg. 4) for this batch, then  $\mathcal{G}_5$  selects the same batch and perturbs its gradient updates with gaussian noise whose magnitude is calibrated to  $C_s$ (Line 6 of Alg. 4). The composition of these two mechanisms uses independent source of randomness through different SGD iterations, hence we can use Theorem 1 to quantify the overall privacy. However, within a single iteration,  $\mathcal{G}_4$  and  $\mathcal{G}_5$  do not use independent source of randomness, although both mechanisms use independent gaussian noise but select the same batch S from the dataset. The following theorem computes  $\alpha_{S_1}(\lambda)$ , and is a generalization of Theorem 1 when the component mechanisms can use dependent source of randomness.

Theorem 3 (General Moments Accountant): Let  $\alpha_{\mathcal{A}_i}(\lambda)$ be  $\max_{D,D'} \log \mathbb{E}_{O \sim \mathcal{A}(D)}[\exp(\lambda \mathcal{P}(\mathcal{A}, D, D', O))]$ , and  $\mathcal{A}_{1:k}$ be the k-fold adaptive composition of  $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k$ . Then:

1) 
$$\alpha_{\mathcal{A}_{1:k}}(\lambda) \leq \sum_{i=1}^{k} j_i \alpha_{\mathcal{A}_i}(\lambda/j_i)$$
  
2)  $\mathcal{A}_{1:k}$  is  $(\varepsilon, \min_{\lambda} \exp(\sum_{i=1}^{k} j_i \cdot \alpha_{\mathcal{A}_i}(\lambda/j_i) - \lambda \varepsilon))$ -DP

for any  $\sum_{i=1}^{k} j_i = 1$ , where  $j_i > 0$  and  $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k$  can use dependent coin tosses. (The proof is in Appendix C.)

Therefore, it follows from Theorem 1 and 3 that:

$$\alpha_{\mathcal{S}_k}(\lambda) \le \alpha_{\mathcal{S}_1}(\lambda) \\ \le T_{\mathcal{S}} \cdot \min_{j_1, j_2 \in (0,1): j_1 + j_2 = 1} \left( j_1 \alpha_{\mathcal{G}_4}(\lambda/j_1) + j_2 \alpha_{\mathcal{G}_5}(\lambda/j_2) \right)$$
(4)

We compute  $\alpha_{\mathcal{G}_4}(\lambda)$  and  $\alpha_{\mathcal{G}_5}(\lambda)$  similarly to [1]. That is, let  $\mu_0(x|\sigma) = g(x|\sigma)$  and  $\mu_1(x|\sigma) = (1-q)g(x|\sigma) + qg(x-1|\sigma)$ , where q = L/|D| is the probability that a record is included in the batch *S* of an SGD iteration and  $g(x|\sigma) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-x^2/2\sigma^2}$ . Then, it holds:

 $\alpha_{\mathcal{G}_3}(\lambda) = \log \max(E_1(\lambda, \sigma_{\mathcal{C}}), E_2(\lambda, \sigma_{\mathcal{C}}))$ 

$$\alpha_{\mathcal{G}_4}(\lambda) = \log \max(E_1(\lambda, \sigma_{\mathcal{S}}), E_2(\lambda, \sigma_{\mathcal{S}}))$$

$$E_1(\lambda,\sigma) = \int_{-\infty}^{\infty} \mu_0(x|\sigma) \cdot \left(\frac{\mu_0(x|\sigma)}{\mu_1(x|\sigma)}\right)^{\lambda} dx$$
$$E_2(\lambda,\sigma) = \int_{-\infty}^{\infty} \mu_1(x|\sigma) \cdot \left(\frac{\mu_1(x|\sigma)}{\mu_0(x|\sigma)}\right)^{\lambda} dx$$

The next theorem immediately follows from Theorem 1 and Theorem 3.

where

<sup>&</sup>lt;sup>4</sup>Computed in the last iteration of Alg. 2.

Dataset	D	$ \mathbb{I}  = m$	$\max   \mathbf{x}  _1$	$\ \mathbf{x}\ _1$
MNIST	60,000	784	311.69	102.44
CDR	4,427,486	1303	422	11.42

TABLE I: Datasets.

*Theorem 4:* DPGM (Alg. 1) is  $(\min_{\lambda} (\alpha_{\mathcal{K}}(\lambda) + \alpha_{\mathcal{S}_k}(\lambda) - \log \delta) / \lambda, \delta)$ -differential private for any fixed  $\delta$ , where  $\alpha_{\mathcal{K}}(\lambda)$  and  $\alpha_{\mathcal{S}_k}(\lambda)$  are defined in Eq. 3 and 4.

In this paper, we use the convention that  $\delta = 1/|D|$ , and compute  $\varepsilon$  numerically. Specifically,  $\varepsilon = \min_{\lambda} (\alpha_{\mathcal{K}}(\lambda) + \alpha_{\mathcal{S}_k}(\lambda) - \log \delta) / \lambda$  is minimized over integer values of  $\lambda$ , where  $\lambda$  is usually no more than 100 in practice. The computation of  $\alpha_{\mathcal{G}_3}$  and  $\alpha_{\mathcal{G}_4}$  are performed through numerical integration, and it suffices to consider 10 different values of  $j_1$  and  $j_2$  in order to have a sufficiently small value of  $j_1 \alpha_{\mathcal{G}_3}(\lambda/\ell_1) + j_2 \alpha_{\mathcal{G}_4}(\lambda/\ell_2)$  in Eq. 4. Therefore, in practice, given  $\delta$ , an accurate approximation of  $\varepsilon$  can be obtained with negligible overhead.

# IV. EXPERIMENTAL EVALUATION

In this section, we compute, numerically, the exact privacy guarantees of DPGM (presented in Alg. 1). Furthermore, we discuss the results of an experimental evaluation assessing its performance in terms of the quality of generated samples as well as counting (linear) queries computed on the synthetic data. Counting queries provide the basis of many data analysis and learning algorithms (see [7] for examples). Finally, we measure the accuracy of our private kernel k-means described in Alg. 2.

Datasets. We use two datasets for our evaluations, summarized in Table I. MNIST is a public image dataset [25], which includes  $28 \times 28$ -pixel images of hand-written digits, a total of 60,000 samples. We vectorize and binarize each image to have binary data records with size m = 784. Throughout our experiments, we assume that each of the 60,000 records originates from a different person. We also use CDR (Call Detail Record) data given to us by a cell phone operator. For this dataset, I represents the set of cell towers of the operator in a large city with |D| = 4,427,486 customers. We use a simplified version of the dataset, which contains the set of visited cell towers per customer within the administrative region of the city over 128.1 km<sup>2</sup>, where the total number of towers is m = 1,303. The average number of individuals per tower over this period was 38,817 with a standard deviation of 50, 911.

**Experimental Settings.** For the RBM, the number of hidden units is set to 200 and the learning rate is 0.01. The biases **b** and **c** are initialized to zeros, while the initial values of the weights **W** are randomly chosen from a zero-mean Gaussian with a standard deviation of 0.01. For the VAE, the number of hidden units is set to 200 with single layer encoder and decoder, and a bi-dimensional latent space. We also used the rectifier activation function (ReLu) for all neurons and the Adam optimizer [23]. For our purposes, it is enough to compute  $\alpha(\lambda)$  for  $\lambda \leq 32$ . We set the number of the private k-means iterations to 20 and  $\delta = 1/|D|$ . We also set  $C_{\text{max}} = 10$ , w = 100 (in Alg. 3), as different values of these parameters do not have a strong impact on the results.

We implement DPGM with both RBM (in C++) and VAE (in Python). Experiments are performed on a workstation running Ubuntu Server 16.04 LTS, with a 3.4 GHz CPU i7-6800K, 32GB RAM, and NVIDIA Titan X GPU card. Source code is available upon request.

**Privacy guarantees.** We report the privacy loss  $\varepsilon$  of DPGM (Alg. 1) in Fig. 2 for the MNIST dataset. Recall that  $\varepsilon$  is computed from the noise level  $\sigma_{\mathcal{C}}$ ,  $\sigma_{\mathcal{K}}$ , and  $\sigma_{\mathcal{G}}$ , the sampling probability q, the number of k-means iterations  $T_{\mathcal{K}}$ , and the number of SGD iterations  $T_{\mathcal{S}}$  using Theorem 4. Fig. 2 shows  $\varepsilon$  depending on the number of SGD training epochs, where one epoch consists of  $\lceil 1/q \rceil$  SGD iterations. In Fig. 2a–2c, we fix  $\sigma_{\mathcal{C}} = 4.0$ , and report the value of  $\varepsilon$  as a function of the number of epochs. We note that larger sampling probabilities (q) and more epochs yield larger values of  $\varepsilon$ , i.e., worse privacy guarantee. Fig. 2b–2c show that larger values of  $\sigma_{\mathcal{K}}$  and  $\sigma_{\mathcal{G}}$  yield stronger privacy guarantees.

Clustering accuracy. Next, in Fig. 4, we compare the private kernel k-means (Alg. 2) with RBF kernel with standard DP kmeans [7]. We evaluate the unsupervised clustering accuracy (ACC) [48], where  $ACC = \max_u \frac{|\{\mathbf{x}: \mathbf{x} \in D \land label(\mathbf{x}) = u(\mathcal{K}(\mathbf{x}))\}|}{|D|}$ ,  $label(\mathbf{x})$  is the ground-truth label of sample  $\mathbf{x}^5$ ,  $\mathcal{K}(\mathbf{x})$  is the cluster assignment obtained by clustering algorithm  $\mathcal{K}$ , and u is a one-to-one mapping between cluster assignments and labels. The best mapping can be obtained using the Hungarian algorithm. To make a fair comparison, we fix  $C_s$  to  $\sqrt{m} = 28$  for standard private k-means without RFF features, and  $C_s = 1$  for private kernel k-means with RFF features based on Theorem 2 - i.e., we do not call DPNorm in either of the algorithms. We compute the clustering accuracy for different values of d depending on  $\sigma_{\mathcal{K}}$ , which directly yields the privacy bound  $\varepsilon$  using Eq. 3 and Theorem 1. Finally, we plot the average accuracy over 100 runs as function of  $\varepsilon$  in Figure 4.<sup>6</sup> Private kernel k-means is clearly superior to standard DP k-means, as the difference in clustering accuracy can be as large as 20%, especially for smaller values of  $\varepsilon$ . Shorter RFF features (i.e., smaller d) result in larger accuracy for smaller values of  $\varepsilon$ , whereas the reverse holds for larger  $\varepsilon$ . For the rest of experiments, we set d to 200. Selecting the optimal number of clusters k for kernel k-means can be qualitatively and visually done by relying on dimensionality reduction algorithms (e.g., t-SNE [?]). To this end, one can use public data sampled from the same underlying distribution, and therefore not requiring to make the parameter selection step differentially private. For MNIST we set k = 10, while we select only one cluster for the CDR dataset. We plan to investigate the effects of different values of k as part of future work.

**Synthetic Samples.** As training progresses, the synthetic samples produced by the generative models should resemble the true samples. To evaluate model quality, we show the synthetic samples obtained at epoch 20 in Fig. 3 from a Restricted Boltzmann Machine and a Variational Autoencoder with k = 10 clusters on MNIST. For this experiment, we set q = 0.0017 for a final privacy budget  $\varepsilon$  of 1.74, and performed  $T_{\mathcal{K}} = 20$  clustering iterations before training the generative neural networks. Overall, the samples generated from VAE

<sup>&</sup>lt;sup>5</sup>For MNIST, these are digits ranging from 0 to 9.

<sup>&</sup>lt;sup>6</sup>Standard deviation of accuracy is less than 0.05 for all values of  $\varepsilon$  and d.



(a) Sampling probability q ( $\sigma_G = 1.0, \sigma_K = 40.0$ ) (b) Clustering noise  $\sigma_K$  ( $\sigma_G = 1.0, q = 0.0017$ ) (c)

(c) SGD noise  $\sigma_{\mathcal{G}}$  ( $\sigma_{\mathcal{K}} = 40.0, q = 0.0017$ )

Fig. 2:  $\varepsilon$  value as a function of the number of SGD training epochs for MNIST ( $\delta = 10^{-5}, T_{\mathcal{K}} = 20$ )



Fig. 3: Real MNIST samples and samples generated from DPGM with RBM and VAE after 20 epochs ( $\varepsilon = 1.74, T_{\mathcal{K}} = 20$ ). In (c) and (d), each row contains 8 samples generated from a cluster.



Fig. 4: Clustering accuracy as a function of  $\varepsilon$  on MNIST ( $\delta = 10^{-5}, T_{\mathcal{K}} = 20$ ).

(Fig. 3c) provide higher visual quality than the ones generated from the RBM (Fig. 3d). Note that the samples generated from the VAE without our private clustering technique (Fig. 3b) have bad visual quality.

**Counting-Query.** We consider counting queries which are specified by a predicate function  $p: D \rightarrow \{0, 1\}$  and return the number of users in the dataset which satisfy the given predicate p, i.e.,  $Q_p(D) = \sum_{\mathbf{x} \in D} p(\mathbf{x})$ . We evaluate the accuracy of counting queries on a synthetic dataset generated by DPGM from our call-data-record (CDR) dataset with roughly 4 million users (see in Table I). A single query is defined by a subset of tower cells, and returns the number of users in D who visited these cells. We compare DPGM with MWEM [21], which is a *de facto* standard differential private mechanism to answer

counting queries. As done in previous work [47], we measure the utility of a counting query  $Q_p$  over the sanitized dataset  $\hat{D}$  by its relative error with respect to the actual result over the raw dataset D. The relative error of  $Q_p$  is thus computed as  $\frac{|Q_p(\hat{D})-Q_p(D)|}{\max\{Q_p(D),s\}}$ , where s is a sanity bound that weight the influence of the queries with small selectivities. Following the convention, the sanity bound is set to 0.1% of the dataset size.

First, we examine the relative error of counting queries with respect to privacy loss  $\varepsilon$ . 1,000 counting queries are randomly generated with different number of tower cells, which we refer as the length of the query. Each query set is divided into 5 subsets such that the query length of the i-th subset is uniformly distributed in  $\left[1, \frac{i \cdot \max ||\mathbf{x}||_1}{\pi}\right]$ and each item is randomly drawn from universe of items. Fig. 5 reports the average relative error for each query set. This shows that our approach clearly outperforms MWEM. The error of DPGM ranges from 0.017 for 20% query length to 0.0012 for 100% when  $\varepsilon = 1.0$ . Weaker privacy guarantee (smaller values of  $\varepsilon$ ) lead to slightly smaller errors (Fig. 5b). By contrast, the error of MWEM<sup>7</sup> ranges from 0.11 to 0.05 even for  $\varepsilon = 2$ . Also note that the synthetic data produced by DPGM allows the evaluation of arbitrary number of type of queries, not only linear counting queries.

<sup>&</sup>lt;sup>7</sup>After clipping each record to have  $L_1$ -norm  $\operatorname{avg}||\mathbf{x}||_1 = 12$ , the sensitivity of queries is set to 12, and the iterations of the algorithm is set to 50 [21].



Fig. 5: Average relative error vs.  $\varepsilon$  ( $q = 2.2 \cdot 10^{-5}, \delta = 4.4 \cdot 10^{-6}$ )

## V. RELATED WORK

In this section, we review prior work on privacy-preserving mechanisms applied to machine learning and data mining.

k-anonymity [40] aims to protect data by generalizing and suppressing certain identifying attributes, however, it does not work well on high-dimensional datasets [3], [9]. Therefore, rather than pursuing input sanitization, prior work has proposed techniques to produce plausible synthetic records with strong privacy guarantees, e.g., focusing on differentially private release of data [2], [11], [13], [22], [29], [30], [45]. Alas, these can often support only the release of succinct data representations, such as histograms or contingency tables. Other mechanisms protect privacy by adding noise directly to the generative model [8], [27], [28], [50]. In this paper, we follow this approach, while, in a first-of-its-kind attempt, focusing on building private generative machine learning models based on neural networks. Other approaches [6], [34], [35] generate data records first, and then attempt to test their privacy guarantees, i.e., decoupling the generative model from the privacy mechanism. By contrast, we attempt to achieve privacy during the training of the model, thus avoiding eventual high sample rejection rates due to privacy tests.

Our work builds on the Differential Privacy (DP) framework, specifically, using the Gaussian mechanism [16]. Due to its generality, DP has served as a building block in several recent efforts at the intersection of privacy and machine learning [1], [38]. In general, the majority of privacy-preserving learning schemes focus on convex optimization problems [4], [12], [46], whereas, training neural networks typically requires to optimize non-convex objective functions - as with Restricted Boltzmann Machine (RBM) [10] and Variational Autoencoder (VAE) [24] – which is usually done through the application of Stochastic Gradient Descent (SGD) with poor theoretical guarantees. Wu et al. [46] propose a private technique which runs SGD for convex cases for a constant number of iterations and only adds noise to the final output. By contrast, in this paper, we introduce a novel differentially private SGD algorithm for optimizing general non-convex loss functions.

Shokri et al. [38] support distributed training of deep learning networks in a privacy-preserving way. Specifically, their system relies on the input of independent entities which aim to collaboratively build a machine learning model without sharing their training data. To this end, they selectively share subsets of noisy model parameters during training. However, their approach incurs high levels of privacy loss per entity, i.e., the  $\varepsilon$ parameter is in the order of thousands. Finally, Abadi et al. [1] introduce an algorithm for non-convex deep learning models with strong differential privacy guarantees. They propose a privacy accounting method, called the moments accountant, which guarantees a tighter bound of the privacy loss for the composition of multiple gaussian mechanisms when compared to the strong composition theorem [17]. Our method also relies on the moments accountant to measure privacy loss, but we train generative models (i.e., unsupervised learning) and with an improved gradient descent, where the noise is carefully adjusted and injected in each iteration.

Differential private k-means has already been addressed in several prior works, for a good overview see [39]. However, all these works aim to find linearly separable clusters, and add noise which is proportional to the data dimension m or the  $L_1$ -norm of data records. By contrast, our private kernel kmeans approach can find even linearly non-separable clusters, and the added noise is independent of d as well as the norm of data points. Also, we offer a tighter privacy analysis using the moments accountant method from [1]. Kernel k-means clustering with random Fourier features (RFF) has already been considered in [14], albeit without any privacy guarantee. Our work somewhat combines [14] and [7], and applies DP kmeans on Fourier features to ultimately achieve better accuracy than [7].

## VI. CONCLUSION

This paper presented a novel differentially private generative model, relying on a mixture of k generative neural networks. The trained models can be used to generate and share synthetic high-dimensional data with provable privacy. We have evaluated the performance of the model on real datasets, showing that our approach provides accurate representation of large datasets with strong privacy guarantees and high utility. In future work, we plan to extend our experiments to location and transportation datasets, as well as pilot deploy our techniques in the wild.

Acknowledgments. Luca Melis was partially supported by The Alan Turing Institute under the EPSRC grant EP/N510129/1. Gergely Acs was supported by the Premium Post Doctorate Research Grant of the Hungarian Academy of Sciences (MTA).

#### REFERENCES

- M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *ACM CCS*, 2016.
- [2] J. M. Abowd and L. Vilhuber. How protective are synthetic data? In PSD, 2008.
- [3] C. C. Aggarwal. On k-anonymity and the curse of dimensionality. In VLDB, 2005.
- [4] R. Bassily, A. Smith, and A. Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In FOCS, 2014.
- [5] Y. Bengio, L. Yao, G. Alain, and P. Vincent. Generalized denoising auto-encoders as generative models. In *NIPS*, 2013.
- [6] V. Bindschaedler, R. Shokri, and C. A. Gunter. Plausible deniability for privacy-preserving data synthesis. *Proceedings of the VLDB Endowment*, 2017.
- [7] A. Blum, C. Dwork, F. McSherry, and K. Nissim. Practical privacy: the SuLQ framework. In *PODS*, 2005.
- [8] C. M. Bowen and F. Liu. Differentially private data synthesis methods. arXiv preprint 1602.01063, 2016.

- [9] J. Brickell and V. Shmatikov. The cost of privacy: destruction of datamining utility in anonymized data publishing. In *KDD*, 2008.
- [10] D. E. Carlson, V. Cevher, and L. Carin. Stochastic spectral descent for restricted boltzmann machines. In AISTATS, 2015.
- [11] A.-S. Charest. How can we analyze differentially-private synthetic datasets? *Journal of Privacy and Confidentiality*, 2011.
- [12] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 2011.
- [13] R. Chen, N. Mohammed, B. C. Fung, B. C. Desai, and L. Xiong. Publishing set-valued data via differential privacy. *Proceedings of the VLDB Endowment*, 2011.
- [14] R. Chitta, R. Jin, and A. K. Jain. Efficient kernel clustering using random fourier features. In *ICDM*, 2012.
- [15] I. S. Dhillon, Y. Guan, and B. Kulis. Kernel k-means: spectral clustering and normalized cuts. In *KDD*, 2004.
- [16] C. Dwork and A. Roth. The Algorithmic Foundations of Differential Privacy. Foundations and Trends in Theoretical Computer Science, 9(3– 4), 2014.
- [17] C. Dwork, G. Rothblum, and S. Vadhan. Boosting and differential privacy. In FOCS, 2010.
- [18] D. Erhan, P. Manzagol, Y. Bengio, S. Bengio, and P. Vincent. The difficulty of training deep architectures and the effect of unsupervised pre-training. In *AISTATS*, 2009.
- [19] M. Fredrikson, S. Jha, and T. Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In ACM CCS, 2015.
- [20] I. Goodfellow, Y. Bengio, and A. Courville. Deep learning. MIT Press, 2016.
- [21] M. Hardt, K. Ligett, and F. McSherry. A simple and practical algorithm for differentially private data release. In *NIPS*, 2012.
- [22] G. Jagannathan and R. N. Wright. Privacy-preserving imputation of missing data. Data & Knowledge Engineering, 2008.
- [23] D. Kingma and J. Ba. Adam: A method for stochastic optimization. arXiv preprint 1412.6980, 2014.
- [24] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *ICLR*, 2014.
- [25] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [26] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. P. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image superresolution using a generative adversarial network. arXiv preprint 1609.04802, 2016.
- [27] H. Li, L. Xiong, and X. Jiang. Differentially private synthesization of multi-dimensional data using copula functions. In *EDBT*, 2014.
- [28] F. Liu. Model-based differential private data synthesis. arXiv preprint 1606.08052, 2016.
- [29] A. Machanavajjhala, D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber. Privacy: Theory meets practice on the map. In *ICDE*, 2008.
- [30] D. McClure and J. P. Reiter. Differential privacy and statistical disclosure risk measures: An investigation with binary synthetic data. *Transactions on Data Privacy*, 2012.
- [31] N. Papernot, M. Abadi, Ú. Erlingsson, I. J. Goodfellow, and K. Talwar. Semi-supervised knowledge transfer for deep learning from private training data. In *ICLR*, 2017.
- [32] J. Pennington, F. X. Yu, and S. Kumar. Spherical random features for polynomial kernels. In *NIPS*, 2015.
- [33] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In NIPS, 2007.
- [34] J. P. Reiter and R. Mitra. Estimating risks of identification disclosure in partially synthetic data. *Journal of Privacy and Confidentiality*, 2009.
- [35] J. P. Reiter, Q. Wang, and B. Zhang. Bayesian estimation of disclosure risks for multiply imputed, synthetic data. *Journal of Privacy and Confidentiality*, 2014.
- [36] T. Salimans, I. J. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *NIPS*, 2016.

- [37] B. Schölkopf, A. J. Smola, and K. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5), 1998.
- [38] R. Shokri and V. Shmatikov. Privacy-preserving deep learning. In ACM CCS, 2015.
- [39] D. Su, J. Cao, N. Li, E. Bertino, and H. Jin. Differentially private k-means clustering. In ACM CODASPY, 2016.
- [40] L. Sweeney. k-anonymity: A model for protecting privacy. International Journal on Uncertainty, Fuzziness and Knowledge-based Systems, 2002.
- [41] L. Theis, W. Shi, A. Cunningham, and F. Huszár. Lossy image compression with compressive autoencoders. arXiv preprint 1703.00395, 2017.
- [42] L. Theis, A. van den Oord, and M. Bethge. A note on the evaluation of generative models. arXiv preprint 1511.01844, 2015.
- [43] T. Tieleman. Training restricted boltzmann machines using approximations to the likelihood gradient. In *ICML*, 2008.
- [44] S. Vempati, A. Vedaldi, A. Zisserman, and C. V. Jawahar. Generalized RBF feature maps for efficient detection. In *British Machine Vision Conference*, 2010.
- [45] L. Wasserman and S. Zhou. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 2010.
- [46] X. Wu, A. Kumar, K. Chaudhuri, S. Jha, and J. F. Naughton. Differentially private stochastic gradient descent for in-RDBMS analytics. arXiv preprint 1606.04722, 2016.
- [47] X. Xiao, G. Bender, M. Hay, and J. Gehrke. ireduct: differential privacy with reduced relative errors. In ACM SIGMOD, 2011.
- [48] J. Xie, R. B. Girshick, and A. Farhadi. Unsupervised deep embedding for clustering analysis. In *ICML*, 2016.
- [49] R. Yeh, C. Chen, T. Lim, M. Hasegawa-Johnson, and M. N. Do. Semantic image inpainting with perceptual and contextual losses. arXiv preprint 1607.07539, 2016.
- [50] J. Zhang, G. Cormode, C. M. Procopiuc, D. Srivastava, and X. Xiao. Privbayes: Private data release via bayesian networks. In *ICMD*, 2014.

#### Appendix

## A. Proof of Lemma 1

Let  $f : \mathcal{D} \to \mathbb{R}$  be a scalar function,  $f(D) = f(D') + \Delta_1 f$ , where  $\Delta_1 f = \Delta_2 f$ , and O = f(D) + x, where  $x \sim \mathcal{N}(0, \sigma)$ . Let  $\hat{\sigma} = \Delta_1 f \cdot \sigma$  Then, it holds:

$$\mathcal{P}(\mathcal{A}, D, D', O) = \ln\left(\frac{\Pr[\mathcal{G}(D) = O]}{\Pr[\mathcal{G}(D') = O]}\right) = \\ = \ln\left(\frac{\Pr[f(D) + \mathcal{N}(0, \hat{\sigma}) = O]}{\Pr[f(D') + \mathcal{N}(0, \hat{\sigma}) = O]}\right) = \ln\left(\frac{\exp(-x^2/2\hat{\sigma}^2)}{\exp(-(x + \Delta_1 f)^2/2\hat{\sigma}^2)}\right) = \\ = \ln\left(\frac{\exp(-x^2/2\hat{\sigma}^2)}{\exp(-(x + \Delta_1 f)^2/2\hat{\sigma}^2)}\right) = \left(\frac{\Delta_1 f}{\hat{\sigma}} \cdot \frac{x}{\hat{\sigma}}\right) + \frac{1}{2}\left(\frac{\Delta_1 f}{\hat{\sigma}}\right)^2$$

Since x is drawn from  $\mathcal{N}(0,\hat{\sigma})$ ,  $\mathcal{P}(\mathcal{A}, D, D', O)$  follows a normal distribution with mean  $(\Delta_1 f)^2/2\hat{\sigma}^2$  and standard deviation  $\Delta_1 f/\hat{\sigma}$ , whose moment generating function is  $\exp\left((\lambda^2 + \lambda)(\Delta_1 f)^2/4\hat{\sigma}^2\right)$ . The claim follows from the definition of  $\alpha$  and  $\hat{\sigma}$ . For the high-dimensional case when  $f: \mathcal{D} \to \mathbb{R}^d$ , the proof is similar to that of Theorem A.1 in [16].

## B. Proof of Theorem 2

Lemma 2: Let  $\mathcal{N}(0,\sigma)$  be a zero-centered normal random variable with standard deviation  $\sigma$ . Then

- 1)  $\mathbb{E}[\cos(\mathcal{N}(0,\sigma))] = \exp(-\sigma^2/2)$  and  $\mathbb{E}[\sin(\mathcal{N}(0,\sigma))] = 0$ ,
- 2)  $\mathbb{E}[\cos^2(\mathcal{N}(0,\sigma))] = (1 + \exp(-2\sigma^2))/2$  and  $\mathbb{E}[\sin^2(\mathcal{N}(0,\sigma))] = (1 \exp(-2\sigma^2))/2$

*Proof:* Let  $\exp(j\mathcal{N}(0,\sigma))$  denote a complex random variable. It follows from the moment generating function of  $\mathcal{N}(0,\sigma)$  that:

$$\mathbb{E}[\exp(j\mathcal{N}(0,\sigma))] = \exp((j\sigma)^2/2) = \exp(-\sigma^2/2)$$

which means that:

$$\mathbb{E}[\cos(\mathcal{N}(0,\sigma)) + j\sin(\mathcal{N}(0,\sigma))] = \mathbb{E}[\exp(j\mathcal{N}(0,\sigma))] = \exp(-\sigma^2/2)$$

This implies that  $\mathbb{E}[\cos(\mathcal{N}(0,\sigma))] = \exp(-\sigma^2/2)$  and  $\mathbb{E}[\sin(\mathcal{N}(0,\sigma))] = 0$  due to the linearity of expectation. Hence

$$\mathbb{E}[\cos^2(\mathcal{N}(0,\sigma))] = \mathbb{E}[(1+\cos(2\mathcal{N}(0,\sigma)))/2] = (1+\exp(-2\sigma^2))/2$$

and

$$\mathbb{E}[\sin^2(\mathcal{N}(0,\sigma))] = \mathbb{E}[(1+\cos(2\mathcal{N}(0,\sigma)))/2] = (1-\exp(-2\sigma^2))/2$$

where we used that  $2\mathcal{N}(0,\sigma) = \mathcal{N}(0,2\sigma)$ .

Proof of Theorem 2: If  $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-\gamma ||\mathbf{x} - \mathbf{y}||_2)$ , then  $p(\mathbf{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^m} \exp(-j \langle \mathbf{w}, \mathbf{x} \rangle) \kappa(\mathbf{w}) d\mathbf{x}$  has zero centered gaussian distribution with standard deviation  $2\gamma \mathbf{I}$ .

$$\mathbb{E}[||z(\mathbf{x})||_{2}] = \mathbb{E}\left[\left((2/d)\sum_{i=1}^{d}\cos^{2}(\langle \mathcal{N}(0,2\gamma\mathbf{I}),\mathbf{x}\rangle + \mathcal{U}[0,2\pi])\right)^{\frac{1}{2}}\right]$$
$$\leq \sqrt{\frac{2}{d}}\left(\sum_{i=1}^{d}\mathbb{E}\left[\cos^{2}(\langle \mathcal{N}(0,2\gamma\mathbf{I}),\mathbf{x}\rangle + \mathcal{U}[0,2\pi])\right]\right)^{\frac{1}{2}}$$
(by Jensen's inequality and the linearity of expectation)

$$\leq \sqrt{\frac{2}{d}} \left( \sum_{i=1}^{d} \mathbb{E} \left[ \cos^{2}(\langle \mathcal{N}(0, 2\gamma \mathbf{I}), \mathbf{x} \rangle)/2 + \sin^{2}(\langle \mathcal{N}(0, 2\gamma \mathbf{I}), \mathbf{x} \rangle)/2 \right] \right)^{\frac{1}{2}}$$
  
$$\leq \sqrt{\frac{1}{d}} \left( \sum_{i=1}^{d} \mathbb{E} \left[ \cos^{2}(\mathcal{N}(0, 2\gamma \sqrt{||\mathbf{x}||_{1}}) \right] + \mathbb{E} \left[ \sin^{2}(\mathcal{N}(0, 2\gamma \sqrt{||\mathbf{x}||_{1}}) \right] \right)^{\frac{1}{2}}$$
  
(by Lemma 2)  
$$\leq 1$$

where, in the second inequality, we used that  $\cos^2(a + b) = \cos^2(a)\cos^2(b) - 2\cos(a)\sin(a)\cos(b)\sin(b) + \sin^2(a)\sin^2(b)$ ,  $\mathbb{E}[\cos(\mathcal{U}[0, 2\pi])] = \mathbb{E}[\sin(\mathcal{U}[0, 2\pi])] = 0$ ,  $\mathbb{E}[\cos^2(\mathcal{U}[0, 2\pi])] = \mathbb{E}[\sin^2(\mathcal{U}[0, 2\pi])] = 0.5$ .

## C. Proof of Theorem 3

This proof is adapted from that of Theorem 2 in [1]. Here we detail the complete proof for the sake of clarity.

Let  $\mathcal{A}_{1:k}$  denote the composition of  $\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_k$  and  $O = (O_1, O_2, \ldots, O_k)$ . Recall that, from Definition 2,  $\mathcal{A}_{1:k}$  is  $(\varepsilon, \delta)$ -DP, if  $\Pr_{O \sim \mathcal{A}_{1:k}(D)}[\mathcal{P}(\mathcal{A}_{1:k}, D, D', O) > \varepsilon] \leq \delta$ .

$$\mathcal{P}(\mathcal{A}_{1:k}, D, D', O) = \log \frac{\Pr[\mathcal{A}_{1:k}(D) = O]}{\Pr[\mathcal{A}_{1:k}(D') = O]}$$

$$= \log \prod_{i=1}^{k} \frac{\Pr[\mathcal{A}_{i}(D) = O_{i}|\mathcal{A}_{i-1}(D) = O_{i-1}, \dots, \mathcal{A}_{1}(D) = O_{1}]}{\Pr[\mathcal{A}_{i}(D') = O_{i}|\mathcal{A}_{i-1}(D') = O_{i-1}, \dots, \mathcal{A}_{1}(D') = O_{1}]}$$
(by the Chain rule)
$$= \sum_{i=1}^{k} \log \frac{\Pr[\mathcal{A}_{i}(D) = O_{i}|\mathcal{A}_{i-1}(D) = O_{i-1}, \dots, \mathcal{A}_{1}(D) = O_{1}]}{\Pr[\mathcal{A}_{i}(D') = O_{i}|\mathcal{A}_{i-1}(D') = O_{i-1}, \dots, \mathcal{A}_{1}(D') = O_{1}]}$$

$$= \sum_{i=1}^{k} \mathcal{P}(\mathcal{A}_{i}, D, D', O_{i})$$
(5)

for any neighboring datasets D and D'. Hence,

 $\alpha_{\mathcal{A}}$ 

$$\leq \sum_{i=1}^{k} j_i \alpha_{\mathcal{A}_i}(\lambda/j_i) \tag{6}$$

where we can apply the generalization of Hölder's inequality in the first inequality due to the fact that  $\exp(\cdot)$  is always positive. Therefore,

$$\left(\sum_{i=1}^{k} \left(\sum_{j=1}^{k} \left(\sum_{i=1}^{k} \left(\sum_{j=1}^{k} \left(\sum_{j=1}^{k}$$

$$\leq \exp(\alpha_{\mathcal{A}_{1:k}}(\lambda) - \lambda\varepsilon) \leq \exp\left(\sum_{i=1} j_i \alpha_{\mathcal{A}_i}(\lambda/j_i) - \lambda\varepsilon\right) \quad \text{(by Eq. (6))}$$

The claim follows from Definition 2.