

M²M: A general method to perform various data analysis tasks from a differentially private sketch

Florimond Houssiau*¹, Vincent Schellekens*², Antoine Chatalic³, Shreyas Kumar Annamraju⁴, and Yves-Alexandre de Montjoye⁴

¹ The Alan Turing Institute

² UCLouvain

³ MaLGa - DIBRIS, Università di Genova, Genoa, Italy

⁴ Imperial College London

*These authors contributed equally.

Abstract. Differential privacy is the standard privacy definition for performing analyses over sensitive data. Yet, its privacy budget bounds the number of tasks an analyst can perform with reasonable accuracy, which makes it challenging to deploy in practice. This can be alleviated by private sketching, where the dataset is compressed into a single noisy sketch vector which can be shared with the analysts and used to perform arbitrarily many analyses. However, the algorithms to perform specific tasks from sketches must be developed on a case-by-case basis, which is a major impediment to their use. In this paper, we introduce the generic *moment-to-moment* (M²M) method to perform a wide range of data exploration tasks from a single private sketch. Among other things, this method can be used to estimate empirical moments of attributes, the covariance matrix, counting queries (including histograms), and regression models. Our method treats the sketching mechanism as a black-box operation, and can thus be applied to a wide variety of sketches from the literature, widening their ranges of applications without further engineering or privacy loss, and removing some of the technical barriers to the wider adoption of sketches for data exploration under differential privacy. We validate our method with data exploration tasks on artificial and real-world data, and show that it can be used to reliably estimate statistics and train classification models from private sketches.

Keywords: Privacy · Differential privacy · Sketching · Sketched learning

1 Introduction

The amount and level of detail of data collected has increased exponentially over the last two decades. Behavioral data has evolved from hand-collected medical records to GPS traces automatically recorded with a temporal resolution on the scale of seconds. While this increased availability and precision of data has resulted in tremendous advances in research, they raise serious privacy concerns. Modern datasets often contain highly detailed summaries of our lives, and are notoriously hard to anonymize. Individuals have indeed been shown to be easily

re-identifiable in large-scale behavioral datasets, such as mobile phone meta-data [27], credit card data [28] and web browsing data [5].

Differential privacy (DP) [14] was introduced by Dwork et al. as a property of algorithms that protect the privacy of users in a dataset. It requires for a randomized algorithm’s outputs to be distributed approximately identically whether any one individual is in the dataset or not. The discrepancy between distributions is controlled by a parameter ε known as the privacy budget. DP is considered by many to as the gold standard definition for privacy loss in aggregated data releases. DP mechanisms have been deployed by institutions with access to large datasets, such as Google to measure changes in mobility patterns caused by confinement measures [2], LinkedIn to answer analytics queries [21], and the US Census for the 2020 Census [1].

Most applications of DP remain limited to specialized tasks on large datasets. Indeed, each differentially private access to a dataset consumes some privacy budget ε , and the total acceptable budget is fixed by the data owner for the dataset. Once this budget has been used entirely, the dataset must be discarded. As such, the number of accurate statistical tasks an analyst can run on a dataset is capped. This strongly limits the utility of differential privacy in practice. In particular, data exploration with DP is particularly challenging: it requires analysts to establish which analyses they want to perform on the dataset, and how to divide the budget between them, before accessing the data.

An increasingly popular solution to this issue is to first compute a differentially private summary of the data, called a private *sketch*, which is then shared with analysts. Once computed, the private sketch can be used as much as desired to solve new learning tasks without accessing the data anymore or using additional privacy budget. This follows from the post-processing property of DP. Sketches have long been used as a technique to compress large-scale datasets to reduce the computational load of algorithms. In this work, the sketch of a dataset is defined as the empirical average of some *feature map* function Φ over all records in a dataset D : $z_D = \frac{1}{|D|} \sum_{x_i \in D} \Phi(x_i)$. The choice of feature map controls the specificity of the information contained in the sketch. For example, researchers have proposed sketches based on Random Fourier Features (RFF) [35] and locality-sensitive hashing [11] that approximate kernel density estimates of the empirical distribution. For some specific sketches and tasks, algorithms with strong theoretical guarantees of accuracy have been developed [17].

However, performing arbitrary data analysis tasks from sketches is difficult, as extracting the desired information from a highly compressed representation of the data is challenging. Each specific task and feature map Φ would require a dedicated algorithm designed by experts. For instance, RFF sketches have in practice only been used for a few tasks, such as Gaussian mixture modeling (GMM) [22] or k-means [23]. Developing compressive methods for other data exploration tasks remains an open problem. This is the main obstacle to using sketches for general data analysis.

In this paper, we introduce a **heuristic to learn from dataset sketches** as shown in Figure 1, which we call the moment-to-moment (M^2M) method.

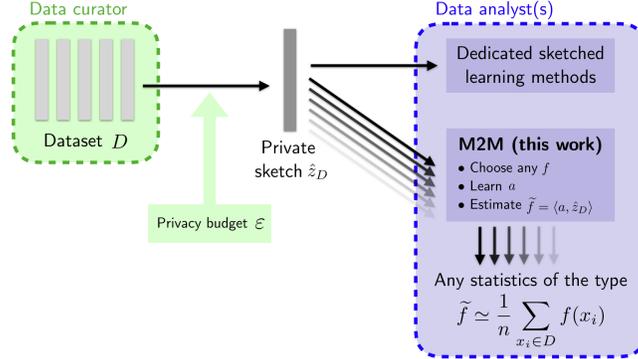


Fig. 1. Considered setup. The data curator releases "once and for all" a private sketch with privacy budget ε . The analyst then chooses a function f and uses our M²M method to learn a vector $a \in \mathbb{R}^m$ such that $\tilde{f} = \langle a, \hat{z}_D \rangle$ approximates the empirical average value of f over the dataset, $\bar{f} = \frac{1}{n} \sum_i f(x_i)$. This procedure can be repeated any number of times (for various choices of f) without additional privacy budget.

M²M allows to approximate empirical averages of functions f from the sketch, $\frac{1}{|D|} \sum_{x_i \in D} f(x_i)$, and can in principle be applied to any feature map Φ . This method is inspired by approximation techniques for kernel methods using random features [32,25]. We **empirically validate our method with artificial and real-world data**, and show that a variety of tasks (moment estimation, counting queries, covariance estimation, logistic regression) can be learned from sketches with comparable performances to alternatives (synthetic data).

2 Background

2.1 Sketches

Sketches are compressed representations of data collections, which can be used to perform some operations efficiently but approximately [12,6]. Sketches usually rely on randomness to achieve a compact representation size. This comes at the price of a probabilistic approximation error. This general principle finds applications in a broad set of contexts, from data streams[16,26] to randomized linear algebra [13]. Here, we focus on sketches that compress the dataset $D = (x_1, \dots, x_n)$ to a single sketch vector $z_D \in \mathbb{R}^m$ by computing the *average of a nonlinear feature map* Φ , applied to each record x_i .

Definition 1. Given a feature map $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$, the sketch of a dataset $D = (x_1, \dots, x_n) \in \mathcal{D}$ is

$$z_D \triangleq \frac{\Sigma_{\Phi}(D)}{|D|} = \frac{1}{n} \sum_{i=1}^n \Phi(x_i) \in \mathbb{R}^m, \quad (1)$$

with $n = |D|$ the dataset size and $\Sigma_{\Phi}(D) = \sum_{i=1}^n \Phi(x_i)$ the sum of features.

The representation $(\Sigma_\Phi(D), |D|)$, where the sum-of-features and dataset size are distinctly encoded, is often used in practice to make it possible to further combine sketches of different datasets into a single one [12].

Typically, sketches are constructed such that scalar products approximate a specific similarity score (called *kernel* $\kappa : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$), $\langle \Phi(x), \Phi(x') \rangle \simeq \kappa(x, x') \forall x, x'$ [30]. This means that they can be used for kernel density estimation (KDE), i.e. building an approximation of the data distribution p_X by a density $\hat{p}(x) \triangleq \frac{1}{n} \sum_{i=1}^n \kappa(x, x_i) \approx \langle \phi(x), z_D \rangle$.

The feature map Φ should be designed such that the sketch z_D captures enough information to solve a target learning task (i.e. the sketched KDE density \tilde{p} accurately represents the true data distribution p_X) while compressing the data as much as possible (i.e. the sketch size m should be small). We here review several important feature map choices.

Histograms Histograms and contingency tables have been extensively studied in the DP literature [14]. Both can be seen as illustrative examples of sketches (in the sense of Definition 1), where the feature map is

$$\Phi^{\text{HIST}}(x) \triangleq (I\{x \in \mathcal{P}_i\})_{i=1, \dots, m} \in \{0, 1\}^m,$$

where $(\mathcal{P}_i)_{i=1}^m$ is a list of subsets of the data domain \mathbb{R}^d , and $I\{A\}$ is the indicator function which returns 1 (resp. 0) whenever A is true (resp. false). For 1-D histograms with n_{bins} bins for example (what we call the HIST sketch), these sets are the one-dimensional bins along each component. For this sketch, $m = d \cdot n_{\text{bins}}$.

RFF Sketches Random Fourier Features (RFF) aim to approximate shift-invariant kernels $\kappa(x, x') = K(x - x')$. They were initially introduced to accelerate kernel methods in machine learning [32].

Definition 2 (Random Fourier Features). *Given $m' = \frac{m}{2}$ “frequency vectors” $\Omega = [\omega_1, \dots, \omega_{m'}] \in \mathbb{R}^{d \times m'}$ drawn $\omega_j \sim_{i.i.d.} \Lambda$, the RFF map is defined as:*

$$\Phi^{\text{RFF}}(x) \triangleq [\cos(x^T \Omega), \sin(x^T \Omega)]^T \in \mathbb{R}^m.$$

The idea is that shift-invariant kernels can be decomposed as $K(x - x') = \mathbb{E}_{\omega \sim \Lambda} e^{i\omega^T(x - x')}$ where the probability distribution Λ is the kernel Fourier transform $\Lambda(\omega) = \int K(u) e^{-i\omega^T u} du$ (owing to Bochner’s theorem [34]). For example, the Gaussian kernel $\kappa(x, x') = \exp(-\|x - x'\|_2^2 / 2\sigma^2)$ admits a Gaussian distribution as Fourier transform, $\Lambda = \mathcal{N}(0, \sigma^{-2} I_d)$. One can then show [32] that up to a constant scaling, Φ^{RFF} satisfies the kernel equation for this kernel.

RFF sketches have been successfully used for parametric density estimation tasks, such as k -means [23] and Gaussian Mixture Modeling [22], reducing the computational resources required by orders of magnitude on large-scale datasets.

RACE Sketches The Repeated Array-of Counts Estimator (RACE) sketch was proposed [11] as an alternative way to approximate KDE for so-called LSH kernels. In RACE, the feature map Φ takes binary values, and is constructed by concatenating R independent hashing functions that each map to W distinct buckets. The size of the sketch is thus $m = R \cdot W$. RACE sketches use *locally-sensitive* hash (LSH) functions: let $W \in \mathbb{N}_0$, a family \mathcal{H} of hash functions $h : \mathbb{R}^d \rightarrow \{1, \dots, W\}$ is locally-sensitive with collision probability κ if $\mathbb{P}_{\mathcal{H}}[h(x) = h(x')] = \kappa(x, x')$ for all $x, x' \in \mathbb{R}^d$.

Definition 3 (Repeated Array-of Counts Estimator). *Given $W \in \mathbb{N}_0$, $h_j, j = 1, \dots, R$, a set of $R = \frac{m}{W}$ hash functions drawn independently from \mathcal{H} , the associated RACE map is defined as:*

$$\Phi^{\text{RACE}}(x) \triangleq [\iota(h_1(x))^T, \dots, \iota(h_R(x))^T]^T \in \{0, 1\}^m,$$

where $\iota : \{1, \dots, W\} \rightarrow \{0, 1\}^W$ denotes the one-hot encoding operation.

Similarly to RFF, one can show [11] that for all choices of LSH, there exists a kernel κ such that the kernel equation is satisfied.

2.2 Differential Privacy

Differential privacy (DP) [14] is seen as the standard definition of privacy for aggregate data releases. It states that the distribution of a differentially private algorithm’s output is similar for any two *neighboring* datasets. Different relations can be considered, but in general (and for the rest of this manuscript), we consider that two datasets are neighbors if they differ by the addition or removal of any one record; this is known as “unbounded” DP⁵. The guarantees of DP are characterized by a privacy “budget” $\epsilon > 0$ which bounds the information disclosure from the dataset. Denote by \mathcal{D} the set of all datasets, equipped with a neighboring relation \sim . In this work, we consider datasets as collections of d -dimensional real-valued vectors $x_i \in \mathbb{R}^d$.

Definition 4 (Differential Privacy). *A randomized mechanism $\mathcal{M} : \mathcal{D} \rightarrow \mathbb{R}^m$ is ϵ -differentially private iff $\forall D \sim D' \in \mathcal{D}, \forall S \subset \mathbb{R}^m$:*

$$\mathbb{P}[\mathcal{M}(D) \in S] \leq e^\epsilon \mathbb{P}[\mathcal{M}(D') \in S].$$

Differential privacy has several desirable properties. First, *composition* guarantees that accessing the same dataset with N different mechanisms respectively using budgets $\epsilon_1, \dots, \epsilon_N$ uses a total budget of at most $\epsilon_{total} = \sum_{i=1}^N \epsilon_i$. Second, *post-processing* ensures that once some quantities have been computed by a differentially private algorithm, no further operation on these quantities can

⁵ We consider only unbounded DP for conciseness, yet the private sketches from Section 2.3 can be extended in a straightforward manner to the bounded DP setting. In this case no noise needs to be added to the denominator in (2).

weaken the privacy guarantees. The latter is particularly important for sketches, as it implies that all analyses ran on a ε -DP sketch are ε -differentially private.

A common method to compute a function f over a dataset with ε -DP is the Laplace mechanism [15]. For a target function $f : \mathcal{D} \rightarrow \mathbb{R}^m$, this mechanism adds centered Laplace noise with scale proportional to the *sensitivity* of f .

Definition 5 (Laplace Mechanism). *The Laplace mechanism to estimate privately a function $f : \mathcal{D} \rightarrow \mathbb{R}^m$ is defined as $\mathcal{M}_f^\zeta(D) = f(D) + \xi$, where $\xi_j \sim \mathcal{L}(\beta)$, $j = 1, \dots, m$ is centered Laplace noise with scale parameter $\beta = \frac{\Delta_1(f)}{\varepsilon}$. The sensitivity $\Delta_1(f)$ is defined as $\Delta_1(f) \triangleq \sup_{D \sim D'} \|f(D) - f(D')\|_1$.*

2.3 Differentially private sketching

We now consider privatized versions of the sketches in the form (1). As the considered feature maps are bounded, their sensitivities are also easily bounded; thus, the Laplace mechanism can be used to produce private versions of these sketches. Following [8] we compute a sketch of the form

$$\hat{z}_D \triangleq \frac{\Sigma_\Phi(D) + \xi}{|D| + \zeta} \triangleq \frac{\sum_{i=1}^n \Phi(x_i) + \xi}{n + \zeta}, \quad (2)$$

where ξ_j ($j = 1, \dots, m$) and ζ are all Laplace random variables with scale parameter chosen according to Definition 5. For ξ , the scale depends on the sensitivity of the sum-of-features function, which can be expressed as $\Delta_1(\Sigma_\Phi) = \max_x \|\Phi(x)\|_1$, which can be computed as: $m'\sqrt{2}$ for RFF [8], R for RACE [11], and k for HIST [14]. For ζ the scale parameter depends on the sensitivity of the cardinality function which is always $\Delta_1(|\cdot|) = 1$. The total privacy budget ε is split across the numerator and the denominator, i.e. the noises ξ and ζ are also respectively proportional to ε_{num}^{-1} and ε_{den}^{-1} , with $\varepsilon = \varepsilon_{num} + \varepsilon_{den}$. As stated above, such private sketches have already been considered in the literature and are not a contribution of this paper: we simply use sketches of this form in order to apply the M²M method introduced in the next section.

Although we focus on pure ε -DP in this manuscript for simplicity, private sketches can easily be extended to satisfy approximate DP (also known as (ε, δ) -differential privacy) using the Gaussian mechanism [15]. This requires computing the L_2 sensitivity of the feature map, see for example [8,18] for RFF.

2.4 Related work

The key advantages of sketching methods for data analysis with differential privacy is that they produce a private “summary” of the dataset, from which an arbitrary number of analyses can be performed. This idea of publishing a DP summary of the data has been explored in the literature, e.g., by Barak et al. with the release of full contingency tables [4]. As contingency tables do not scale with the number of dimensions, further work has been proposed to publish so-called “views” of the data, from which n -way marginals can then be computed [31]. Another type of data summary that has gained popularity in recent years is

synthetic data, where the data curator publishes a dataset is “similar” to the original data, but with no mapping from real to synthetic records. These usually involve training a statistical model on the data, which is then used to generate synthetic records, either explicitly [24,37] or using generative networks [36].

Kernel mean embeddings are known to carry a lot of information on the data distribution and are thus of particular interest for privacy applications. Balog et al. suggested to use synthetic data points in order to represent (possibly infinite-dimensional) kernel mean embeddings in a private manner [3]. Finite-dimensional approximations based on random Fourier features have been made private using simple additive perturbation mechanisms with applications to clustering and Gaussian modeling [8] as well as synthetic data generation [18]. More recently, compact sketches based on Hermite polynomials have been proposed [29] and have been shown empirically to provide a better privacy-utility tradeoff for private data generation than random Fourier features.

Relating specifically to the M²M method, the idea of considering a learned linear combination of random features (without privacy) has been popularized by Rahimi and Recht [33], and then extensively studied under the name of “extreme learning machines” (ELMs) [20,19]. The sketches considered in this paper can be interpreted as instances of this idea, with an additional averaging operation over the dataset. This is made possible by the fact that we only consider learning moments of the data.

3 The moment-to-moment method

Sketching methods can be used to efficiently perform specific learning tasks, and can often be made private in a straightforward manner by additive perturbation; however, extracting information from them is hard in general. Here, we introduce the *moment-to-moment* (M²M) heuristic to learn a broad range of aggregate statistics from a single sketch. While previous sketched learning methods were relatively specific in the sense that both the feature map and the algorithm to learn from the sketch had to be tailored to a specific machine learning task, our heuristic can be used to approximate various kinds of statistics from the same sketch. Although M²M can naturally be used on a non-private sketch, it is particularly attractive for private sketches, as it allows an analyst to perform arbitrarily many analyses from the sketches without incurring any additional privacy budget.

In the following, we assume that the data curator holds a sensitive dataset D of size n , chooses a data-independent feature map Φ , and releases publicly the triplet $(\Phi, \hat{z}_D, n + \zeta)$ where $\hat{z}_D = \frac{1}{n+\zeta}(\sum_{i=1}^n \Phi(x_i) + \xi)$ is the private sketch computed as in (2) and ξ, ζ are random and chosen as explained in Section 2.3 in order to ensure ϵ -DP (i.e. they depend on the sensitivity of the feature map Φ). Note that any result obtained by post-processing from this triplet will always remain ϵ -DP.

3.1 Method description

Suppose that an analyst wants to compute the empirical average of an arbitrary target function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ over the dataset, i.e. the quantity $\bar{f} \triangleq \frac{1}{n} \sum_{i=1}^n f(x_i)$. The M²M method estimates \bar{f} by a linear function of the sketch $\langle a, \hat{z}_D \rangle$, where the coefficients $a \in \mathbb{R}^m$ are computed by the method. Because both the input (the sketch \hat{z}_D) and the output (the target average \bar{f}) can be seen as “generalized moments” (averages of some features of the data) of the dataset, this amounts to transforming one type of generalized moment to another, hence the name of our method.

In order to apply the method, an analyst chooses *a priori* a bounded domain $D_{M^2M} \subset \mathbb{R}^d$ such that all possible records lie inside of D_{M^2M} (for example, D_{M^2M} might be a box constrained by physical upper and lower bounds on the data values). The principle of M²M is to approximate the target function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ over this domain D_{M^2M} by a linear model \tilde{f} of parameters $a \in \mathbb{R}^m$ in the output space of the sketch feature map $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$, i.e.,

$$\tilde{f}(x) \triangleq \langle a, \Phi(x) \rangle \approx f(x), \quad \forall x \in D_{M^2M}.$$

The key insight is that this linear model can then be used to estimate the dataset average \bar{f} from the dataset sketch z_D , since the sketching operator is linear:

$$\tilde{f}(z_D) = \langle a, z_D \rangle = \frac{1}{n} \sum_{i=1}^n \langle a, \Phi(x_i) \rangle \approx \frac{1}{n} \sum_{i=1}^n f(x_i) = \bar{f}. \quad (3)$$

Intuitively, the target function $f(\cdot)$ is approximated by a linear combination $\langle a, \Phi(\cdot) \rangle$ of a set of m base functions (the components of the feature map, $\Phi_i(\cdot)$). The quality of the approximation depends on the compatibility between the feature map Φ and the function to approximate f . Zhang et al. [38] showed that such a linear combination can approximate continuous functions arbitrary well for a large enough number of features m , under conditions satisfied by many standard feature maps. This suggests that M²M can be used to approximate any continuous function f for a large array of sketches, although quantifying precisely how the approximation quality decreases with m is out of the scope of this paper. It should also be expected that approximating a discontinuous function f with, e.g., RFF features will lead to high approximation error (e.g., some kind of Gibbs phenomenon).

We illustrate M²M with a toy example in Fig. 2. We consider the step function $f(x) = I\{x \geq 0.5\}$ restricted to the domain $D_{M^2M} = [0, 1] \subset \mathbb{R}^{d=1}$. For RFF, the approximation $\tilde{f}(x)$ is a linear combination of $\cos(\omega_i^T x)$ and $\sin(\omega_i^T x)$, for some fixed ω_i , which explains the bumps observed in the approximation (Gibbs phenomenon). The RACE feature map, whose base functions are the one-hot encoding of locally-sensitive hash functions, approximates f by a piecewise constant function.

3.2 Optimizing the M²M model

For the results produced by the M²M method to be useful, the parameters of the linear model a need to be chosen such that $\tilde{f}(\cdot)$ is a good approximation of the

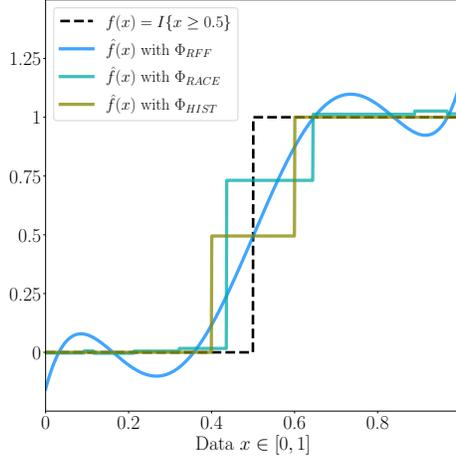


Fig. 2. How M²M works: the M²M method approximates the target function f as a linear combination of components of the feature map, $\sum_{i=1}^m a_i \Phi_i(x) \approx f(x)$.

true function $f(\cdot)$ on the domain of interest. For this, we formulate and optimize a loss function J for the vector of weights a that penalizes differences between f and \tilde{f} . The full learning procedure is described in algorithm 1 in appendix B.

Similarly to Rahimi et al. [32], we use the squared difference as distance, $d(f(x), \tilde{f}(x)) = (\tilde{f}(x) - f(x))^2$. Assume that the records $x_i \in D$ are drawn from some (unknown) probability distribution $x_i \sim_{i.i.d.} p_X$. Ideally, the M²M procedure would minimize the average error of the approximation *over the true data distribution* p_X , $J_{\text{ideal}}(a) = \mathbb{E}_{X \sim p_X} \left[d(f(x), \tilde{f}(x))^2 \right]$. However, the analyst only has access to the private sketch and does not know p_X , let alone the data D . Instead, they choose an *a priori* distribution ψ that is either (1) close to p_X , or (2) likely to yield a good approximation where p_X takes significant values when optimizing for the (approximate) loss $J_\psi(a) = \mathbb{E}_{X \sim \psi} \left[d(f(x), \tilde{f}(x))^2 \right]$. In this work, we assume no prior knowledge except for the domain $D_{\text{M}^2\text{M}}$ and thus use the uniform distribution on this domain $\psi = \text{Unif}(D_{\text{M}^2\text{M}})$, following the principle of maximum entropy. Finally, since evaluating the expectation operator analytically can be challenging for arbitrary ψ , f and Φ , especially in high dimensions, we approximate it by sampling a large number n_s of *training synthetic* data points $(\tilde{x}_i)_{i=1}^{n_s}$ sampled i.i.d. from ψ . The resulting loss, given a choice of ψ , is:

$$J_{\text{noreg}}(a) = \frac{1}{N} \sum_{i=1}^N (f(X_i) - \langle a, \Phi(X_i) \rangle)^2 \quad X_1, \dots, X_N \sim_{i.i.d.} \psi$$

However, minimizing J_{noreg} directly is not robust to noise, and in particular to the noise added to obtain differential privacy. Indeed, when applying the linear model a from (3) to the private data summary \hat{z}_D , we get (neglecting, for

illustration, the noise ζ on the denominator):

$$\langle \hat{z}_D, a \rangle = \left\langle \frac{1}{n} \sum_{i=1}^n \Phi(x_i) + \xi, a \right\rangle \approx \frac{1}{n} \sum_{i=1}^n f(x_i) + \frac{1}{n} \langle \xi, a \rangle.$$

Hence, the noise on the numerator ξ causes an error in the M²M estimate of variance $\sigma_\xi^2 \|a\|_2^2 / n^2$. To account for this noise, we add a term proportional to its variance to the loss function J :

$$J(a) \triangleq \mathbb{E}_{X \sim \psi} \left[(f(X) - \langle a, \Phi(X) \rangle)^2 \right] + \lambda \|a\|_2^2, \quad (4)$$

where we set the regularization parameter λ to the value σ_ξ^2 / n^2 . We prove that this loss J is an upper bound for the mean square prediction error between \bar{f} and the M²M estimate $\langle a, \hat{z}_D \rangle$ (see proof in Appendix A).

Theorem 1. *Let $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a feature map, $D_{M^2M} \subset \mathbb{R}^d$, and D be a random dataset of n records $X_1, \dots, X_n \sim_{i.i.d.} \psi$. For all $a \in \mathbb{R}^m$, and all distributions ψ , if $\lambda = \sigma_\xi^2 / n^2$ and $\zeta = 0$, we have that, if $\zeta = 0$:*

$$J(a) \geq \mathbb{E}_{X_1, \dots, X_n, \xi} \left[\left(\frac{1}{n} \sum_{i=1}^n f(X_i) - \langle a, \hat{z}_D \rangle \right)^2 \right]$$

Since the exact dataset size n is not directly available to the analyst, we use $|D| + \zeta$ as an estimation of n . Further to this, we found empirically that using $\lambda = \frac{\sigma_\xi^2}{n^2}$ makes the model insufficiently robust to noise (especially when the sensitivity of the feature map is large). We thus use a larger regularization parameter in experiments by removing the square on the estimated number of samples.

$$\lambda = \frac{\sigma_\xi^2}{(|D| + \zeta)} = \frac{2 \cdot \Delta_1(\Phi)^2}{\varepsilon_{num}^2 \cdot (|D| + \zeta)}. \quad (5)$$

Solving for J Let $(\tilde{x}_i)_{i=1}^{n_s}$ denote the set of random training samples used inside the M²M procedure. Denote the synthetic feature matrix $\mathbf{P} = (\Phi(\tilde{x}_i))_{i=1}^{n_s} \in \mathbb{R}^{n_s \times m}$, and the vector of corresponding outputs $\mathbf{F} = (f(\tilde{x}_i))_{i=1}^{n_s} \in \mathbb{R}^{n_s}$. The empirical loss that M²M optimizes is $J(a) = \frac{1}{n_s} \|\mathbf{P} \cdot a - \mathbf{F}\|_2^2 + \lambda \|a\|_2^2$. This corresponds to a ridge regression problem with regularization parameter λ , and can be solved efficiently.

3.3 Sources of error

M²M is a heuristic method to approximate \bar{f} , and as such will always incur some error. We here outline the four main sources of error of M²M.

1. *Sampling error*: The expectation operator in the cost function $J(a)$ is not computed exactly, but estimated by sampling n_s points $\tilde{x}_i \sim \psi$. If n_s is too small, this estimate can be inaccurate, and the model a risks “overfitting” to the small training set.

2. *Approximation error*: M²M finds coefficients a such that the linear combination $\tilde{f}(\cdot) = \langle a, \Phi(\cdot) \rangle = \sum_{i=1}^m a_i \Phi_i(\cdot)$ approximates the target function f . In general, even if a is the exact minimizer of $J(a)$, there remains some inherent approximation error which depends on the compatibility between the feature map Φ and target function f .
3. *Distributional shift*: In practice, the empirical distribution p_X differs from the probability distribution ψ used for training. Distributional shift is a hard problem to fix, as it requires tailoring ψ to p_X *without accessing the data*, or only through the sketch. We discuss this in Section 5.
4. *Differential privacy noise*: Finally, the noises ξ and ζ added in the computation of the sketch \hat{z}_D further distort the representation. This error decreases when the privacy budget ε increases.

3.4 Statistical estimation with M²M

Many learning tasks can be written as the estimation of some generalized moments of the data. Here we give some common examples.

1. Moments: The j^{th} component of the k^{th} moment of the empirical data distribution is defined as

$$m_j^{(k)} = \frac{1}{n} \sum_{i=1}^n (x_i)_j^k \approx \mathbb{E}_{X \sim p_X} X_j^k,$$

which is the empirical average of the function $f^{(j,k)} : \mathbb{R}^d \rightarrow \mathbb{R} : x \mapsto x_j^k$.

2. Counting queries: Given a set $S \subset \mathbb{R}^d$, a counting query over D consists of finding the number of data points from the dataset D that belong to S :

$$\text{COUNT}(D, S) = |\{i \in \{1, \dots, n\} : D_i \in S\}| = \sum_{1 \leq i \leq n} f_S(x_i).$$

where $f_S : \mathbb{R}^d \rightarrow \{0, 1\} : x \mapsto I\{x \in S\}$ denotes the indicator function of S . Histograms are a specific subset of counting queries, where the set S is chosen to be a one-dimensional “bin”.

3. Covariance: The $(i, j)^{\text{th}}$ entry of the empirical covariance matrix is

$$c_{ij} = \frac{1}{n} \sum_{l=1}^n ((x_l)_i - \mu_i) \cdot ((x_l)_j - \mu_j),$$

which is the empirical average of the function $f^{(i,j)} : \mathbb{R}^d \rightarrow \mathbb{R} : x \mapsto (x_i - \mu_i)(x_j - \mu_j)$. The mean of the component i , μ_i , can be estimated using M²M for the first-order moment, $m_i^{(1)}$.

3.5 Classification and regression by approximation of the loss

Many learning tasks can be formulated as learning a parametric model with parameter θ using a loss function L . For such tasks, one will typically solve the optimization problem $\theta^* \in \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n L(x_i, \theta)$, whose objective function takes the form of a generalized moment. Specifically, for a classification or regression task, the analyst wants to fit some model $F_{\theta} : \mathbb{R}^{d-1} \rightarrow \mathbb{R}$ parameterized by $\theta \in \mathbb{R}^p$ to the data samples $(x_i)_{i=1}^n$, where each sample x_i is a pair

$x_i = (\bar{x}_i \in \mathbb{R}^{d-1}, y_i \in \mathbb{R})$. If the fitting quality is quantified by a loss function $l(\cdot, \cdot)$, one can define $L_\theta(x_i) \triangleq L(x_i, \theta) \triangleq l(F_\theta(\bar{x}_i), y_i)$ and M²M can be used with the target $f = L_\theta$ for any *fixed* value of θ . Finding the optimal parameter θ^* involves solving the following bi-level optimization problem:

$$\theta^* \in \arg \min_{\theta} \langle a_\theta, \hat{z}_D \rangle \quad \text{such that} \quad a_\theta \in \arg \min_a J_\theta(a) \quad (6)$$

where J_θ is the M²M objective associated to the target function L_θ . As mentioned in Section 3.2, solving for a is a ridge regression a problem, which has a closed-form solution (given synthetic samples \tilde{x} used to compute J) of $a_\theta = \mathbf{S} \cdot \sum_{i=1}^{n_s} \Phi(\tilde{x}_i) L_\theta(\tilde{x}_i)$ where $\mathbf{S} = \left(\frac{1}{n_s} \sum_{i=1}^{n_s} \Phi(\tilde{x}_i)^T \Phi(\tilde{x}_i) + \lambda I \right)^{-1}$. We then use this result in eq. 6 to formulate the dual optimization problem as an optimization problem in θ^* :

$$\theta^* \in \arg \min_{\theta} \sum_{i=1}^{n_s} \underbrace{\Phi(\tilde{x}_i)^T \cdot \mathbf{S} \cdot \hat{z}_D}_{\triangleq w(\tilde{x}_i)} \cdot L_\theta(\tilde{x}_i)$$

This method, which we call *implicit-M²M*, computes a weighting function $w : \Omega \rightarrow \mathbb{R}$ from the feature map Φ , private sketch \hat{z}_D , and regularization coefficient λ , independently of the loss. This weighting function is used to weigh the contribution of each synthetic points to the total loss. Any learning procedure, such as gradient descent, can then be applied to the re-weighted loss.

4 Experiments

We empirically evaluate the M²M method on a range of data analysis tasks on artificial and real data. We perform our analyses on the LifeSci dataset, a real-world dataset of life sciences measurements ($n = 2.7 \cdot 10^4$ records and $d = 10$ attributes), which we normalize to $\Omega = [0, 1]^d$. In order to analyze the different sources of errors independently, we perform the same analyses on a uniformly sampled artificial dataset of same shape (n, d) , which we call **Random10**. Since the training distribution ψ is equal to the empirical distribution p_X , there is no *distributional shift*, and the error observed in the results for **Random10** is thus the combination of *approximation error*, *sampling error* and *the DP noise addition*.

We sketch each dataset the using RFF ($m = 200, \sigma = 1$), RACE ($R = W = 80$), and HIST (marginals of each attribute, $n_{bins} = 100$ bins of same size in $[0, 1]$), and add noise to ensure DP with privacy budget $\epsilon \in [10^{-2}, 10^2]$, as described in Section 2.3. For all sketches, we split the privacy budget as $\epsilon_{num} = 0.98 \epsilon$ and $\epsilon_{den} = 0.02 \epsilon$. We train M²M models with $n_s = 10^5$ samples, which empirically results in very low sampling error (training and testing R^2 scores almost identical). We repeat each experiment 50 times and report, for each task, the average accuracy over all runs.

An alternative to sketches is synthetic data generation (SDG), where a statistical model is fit to the real data, and so-called synthetic data are then generated

by sampling from this model. We compare our results with datasets generated using three differentially private SDG methods: **DP-Copula** [24], **PrivBayes** [37], and **DP-WGAN** [36]. The latter method relies on a relaxed definition of differential privacy, (ϵ, δ) -DP, and hence the guarantees provided are weaker. In our experiments, we use $\delta = 10^{-5}$. For each SDG and ϵ , we generate 10 synthetic datasets from LifeSci, and perform the tasks of interest on the synthetic data (by computing the empirical average of the functions f on the synthetic data), reporting the average over all runs.

4.1 Tasks involving columns in isolation

As a first illustrative example, we consider a range of simple tasks where the function learned with M²M only concerns one attribute in isolation. For each sketch and each column in the datasets, we train a M²M model to predict (1) its mean $\frac{1}{n} \sum_{i=1}^n x_i$, (2) its order 2 moment $\frac{1}{n} \sum_{i=1}^n x_i^2$, and (3) its cumulative distribution function (CDF) in 10 equi-distant points $(\frac{1}{n} \sum_{i=1}^n I\{x_i \leq S_j\})_{j=1}^{10}$. We then measure the error obtained between the predicted value and the empirical value using mean relative error (MRE) $MRE(\hat{\mu}, \mu) = \frac{|\hat{\mu} - \mu|}{\mu}$ for (1) and (2), and the Earth-Mover Distance (EMD) for (3). For each task, sketch, and dataset, we report the average error across all attributes in Fig. 3.

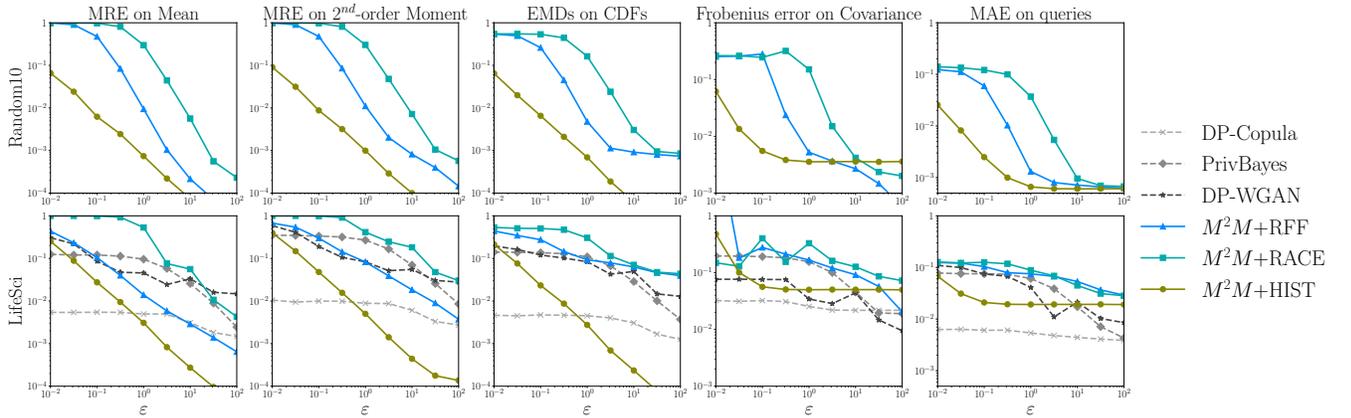


Fig. 3. Estimation of one-dimensional statistics over a random artificial dataset (top row) and LifeSci (bottom row). We estimate the mean, second-order moment, and CDF of each attribute using M²M on three sketches (RFF, RACE, and HIST), and synthetic datasets (generated using DP-Copula, PrivBayes, and DP-WGAN). We estimate the covariance matrix and the answer to a large number of random counting queries using M²M on three sketches (RFF, RACE, and HIST), and synthetic datasets (generated using DP-Copula, PrivBayes, and DP-WGAN).

We show that, in the absence of distributional shift, M²M can be used to estimate single-variable tasks with good accuracy. As expected, the HIST sketch performs well on all tasks and for both datasets, since it is specifically designed to approximate one-dimensional distributions. However, distributional shift (in LifeSci) worsens results significantly for all feature maps. This is particularly true for CDF, where the RFF and RACE feature maps result in high error, probably due to high approximation error. Comparing with synthetic data, we find that the RFF sketch compares favorably with both PrivBayes and DP-WGAN (especially when $\varepsilon \geq 1$), while the RACE sketch leads to less useful results. DP-copula datasets outperform both sketches, which is to be expected since the method explicitly estimates marginals.

We further analyze the different sources of error in table 1. We report the mean relative error on the first moment $\mathbb{E}[X]$ obtained with either the exact sketch ($\varepsilon = +\infty$) or the private sketch with parameter $\varepsilon = 1$, for the RFF and HIST feature maps, on both datasets. For the HIST feature map and $\varepsilon = +\infty$, we find the M²M coefficients using a small regularization $\lambda = 10^{-9}$ (for numeric stability). The error on the artificial dataset for $\varepsilon = +\infty$ is the *approximation error* of f , the irreducible error obtained when approximating f by a linear mixture of components of the feature map Φ_i . We observe that this error is low for the RFF feature map, which has strong approximation properties [32,38], and higher for the HIST sketch, which roughly approximates a function as a product of 1D piecewise constant functions. The second row ($\varepsilon = 1$) is the result of adding *DP error* to the approximation error. DP error has a negligible impact on the performances of the histogram sketch, as it is dominated by the approximation error. The opposite applies to RFF, where the DP error is 5 orders of magnitude larger. Results from the LifeSci dataset (rows 3 and 4) illustrate the impact of distributional shift, when the distribution used to generate M²M’s training set differs from the empirical distribution. For $\varepsilon = 1$, we observe that all resulting errors are one order of magnitude larger, as a result of distributional shift. Furthermore, as expected, when there is no DP error ($\varepsilon = +\infty$), the approximation error for LifeSci is higher than for the Random10, for both sketches. Hence, distributional shift can have disparate effects on the resulting accuracy of the method, by amplifying either or both of the approximation and DP error.

Dataset	Budget	MRE for Φ^{RFF}	MRE for Φ^{HIST}
Random10	$\varepsilon = +\infty$	$6.25 \cdot 10^{-8}$	$1.87 \cdot 10^{-5}$
	$\varepsilon = 1$	$9.55 \cdot 10^{-3}$	$9.10 \cdot 10^{-4}$
LifeSci	$\varepsilon = +\infty$	$1.67 \cdot 10^{-6}$	$5.91 \cdot 10^{-5}$
	$\varepsilon = 1$	$4.20 \cdot 10^{-2}$	$3.8 \cdot 10^{-3}$

Table 1. Comparison of asymptotic, DP, and distributional shift errors: We measure the RMSE on the first moment $\mathbb{E}[X]$ estimated with the M²M method and the Random Fourier Features Φ^{RFF} and HIST Φ^{HIST} feature maps, on the artificial and LifeSci datasets. We report the asymptotic error (no noise) and the error for $\varepsilon = 1$. All results are averaged over 100 trials.

4.2 Multi-column tasks

We evaluate M²M on tasks that involve attributes taken together. First, we compute the covariance matrix of the dataset, $\frac{1}{n} \sum_{i=1, j=1}^{n, n} (x_i - \hat{\mu}_i) \cdot (x_j - \hat{\mu}_j)$, using $\hat{\mu}_i$ estimated as above. We measure the Frobenius distance between the estimated and empirical covariance matrices. Next, we perform a large number of simple counting queries $\text{COUNT}(D, S)$, where the query S is defined as the conjunction of three predicates of the form $X_i \leq u$ or $X \geq l$, for three different attributes $X_i, X_{i'}, X_{i''}$. We report the Mean Absolute Error (MAE) between the real query answers and the answers predicted by M²M.

Fig. 3 reports the error decrease for both tasks and on each dataset as ϵ increases. Similarly to the one-dimensional tasks (Fig. 3), we observe that M²M estimations perform well on the Random10 dataset, and worse on LifeSci. Except for PrivBayes, all synthetic datasets (and in particular, DP-Copula) outperform M²M. The queries use case is particularly challenging to approximate with M²M, as the target function f is not continuous. Finally, as expected, results for the HIST sketch quickly plateau for all tasks and datasets.

4.3 Logistic Regression

We use the implicit-M²M method described in section 3.5 to perform logistic regression from the private sketch of a dataset. We use real-world building occupancy data [7] ($d = 6, n = 20,560$) with 5 continuous attributes (building characteristics) and a binary attribute (whether a building is occupied). This dataset is such that the last attribute is strongly predicted by the continuous attributes, with an AUC (area under curve) of > 0.99 . We normalize the continuous attributes to $[0, 1]$ and define $\Omega = [0, 1]^5 \times \{0, 1\}$ and $\psi = \text{Unif}(\Omega)$. We randomly separate the data between training (90%) and testing (10%), then sketch the training dataset using RFF ($\sigma = 1, m = 200$), RACE ($R = 80, H = 80, \sigma = 0.1$) and HIST ($n_{bins} = 20$) for a range of ϵ . Using implicit-M²M, we perform logistic regression on each sketch and evaluate the result on the testing dataset. We compare our results with Chaudhuri et al. ’s DP-ERM [9], a dedicated method to train a logistic regression with DP using objective perturbation.

We also generate synthetic datasets using the same SDG techniques as above. We train a logistic regression using `sklearn` on each dataset 10 times, and measure its AUC on the test dataset. It can happen that the synthetic dataset only has one class for the last attribute; in this case we report the AUC to be 0.5.

In Fig. 4, we show that implicit-M²M compares remarkably well with DP-ERM for the RFF feature map. While it leads to higher error, the RACE feature map consistently produces an AUC of at least 0.9 for $\epsilon \geq 0.3$. Unsurprisingly, the method performs poorly on the HIST feature map ($AUC < 0.1$, not featured on the plot), which cannot, by definition, be used to estimate correlations between attributes. Importantly, models trained with implicit-M²M compare favorably with models trained on synthetic datasets using the same budget ϵ . As expected, the task-specific DP-ERM outperforms all other methods, but this comes at the cost of the entire budget ϵ . Our results suggest that implicit-M²M is a promising solution to perform sophisticated learning tasks on sketches.

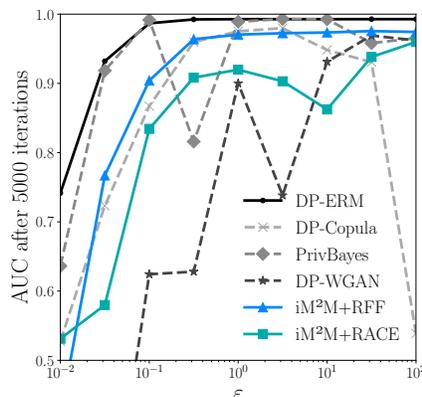


Fig. 4. AUC of a logistic regression trained from sketches on the occupancy dataset. We use implicit-M²M to fit a logistic regression to the occupancy dataset from RFF and RACE sketches. We compare our results with the dedicated method DP-ERM and three synthetic data generation methods.

5 Future work and conclusion

Distributional shift occurs when the distribution used to generate M²M’s training set, ψ , differs from the data distribution p_X . This is a significant source of error in the method. We here propose a few options to reduce this error.

- Improving the approximation $\psi \approx p_X$ using the sketch. KDE sketches are built to approximate a kernel, encoding a kernel density estimate for the data distribution: $p_X(x) \approx \frac{1}{n} \sum_{i=1}^n \kappa(x, x_i) \approx \langle \Phi(x), \hat{z}_D \rangle$. One could thus use $\psi : x \mapsto \langle \Phi(x), \hat{z}_D \rangle$. However, the approximate distribution $\langle \Phi(x), \hat{z}_D \rangle$ can be negative and is not robust to noise addition for privacy.
- Learning a generative model on the sketch [18] that, if accurately trained, generates synthetic data similar to the sketched dataset. These synthetic records can then be used to train the M²M model, as their distribution p_{synth} is likely to be close to p_X (or at least closer than ψ uniform). Although the synthetic records could be used directly for the learning tasks, re-accessing the data sketch through the M²M mechanism could yield greater utility.
- Solving the loss minimization problem on the real data using a differentially private procedure. For instance, techniques such as DP-Empirical Risk Minimisation (DP-ERM) [10] could be applied – although this can be challenging, since J is non-convex. While this method is most likely the best solution to distributional shift, it requires additional privacy budget to learn the parameters of M²M, which contradicts the idea of data summaries.

References

1. Abowd, J.M.: The us census bureau adopts differential privacy. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. pp. 2867–2867 (2018)
2. Aktay, A., Bavadekar, S., Cossoul, G., Davis, J., Desfontaines, D., Fabrikant, A., Gabrilovich, E., Gadepalli, K., Gipson, B., Guevara, M., et al.: Google covid-19 community mobility reports: Anonymization process description (version 1.0). arXiv preprint arXiv:2004.04145 (2020)
3. Balog, M., Tolstikhin, I., Schölkopf, B.: Differentially Private Database Release via Kernel Mean Embeddings. In: International Conference on Machine Learning. pp. 414–422 (Jul 2018)
4. Barak, B., Chaudhuri, K., Dwork, C., Kale, S., McSherry, F., Talwar, K.: Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In: Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems. pp. 273–282 (2007)
5. Barbaro, M., Zeller, T., Hansell, S.: A face is exposed for aol searcher no. 4417749. *New York Times* **9**(2008), 8For (2006)
6. Blum, A., Hopcroft, J., Kannan, R.: Foundations of Data Science. Cambridge University Press (2020)
7. Candanedo, L.M., Feldheim, V.: Accurate occupancy detection of an office room from light, temperature, humidity and co2 measurements using statistical learning models. *Energy and Buildings* **112**, 28–39 (2016)
8. Chatalic, A., Schellekens, V., Houssiau, F., De Montjoye, Y.A., Jacques, L., Gribonval, R.: Compressive learning with privacy guarantees. *Information and Inference: A Journal of the IMA* (iaab005) (May 2021). <https://doi.org/10.1093/imaiai/iaab005>
9. Chaudhuri, K., Monteleoni, C.: Privacy-preserving logistic regression. In: Advances in neural information processing systems. pp. 289–296 (2009)
10. Chaudhuri, K., Monteleoni, C., Sarwate, A.D.: Differentially private empirical risk minimization. *Journal of Machine Learning Research* **12**(3) (2011)
11. Coleman, B., Shrivastava, A.: Sub-linear race sketches for approximate kernel density estimation on streaming data. In: Proceedings of The Web Conference 2020. pp. 1739–1749 (2020)
12. Cormode, G., Garofalakis, M., Haas, P.J., Jermaine, C.: Synopses for massive data: Samples, histograms, wavelets, sketches. *Foundations and Trends in Databases* **4**(1–3), 1–294 (2012)
13. Drineas, P., Kannan, R., Mahoney, M.W.: Fast monte carlo algorithms for matrices i: Approximating matrix multiplication. *SIAM Journal on Computing* **36**(1), 132–157 (2006)
14. Dwork, C.: Differential privacy: A survey of results. In: International conference on theory and applications of models of computation. pp. 1–19. Springer (2008)
15. Dwork, C., McSherry, F., Nissim, K., Smith, A.: Calibrating noise to sensitivity in private data analysis. In: Theory of cryptography conference. pp. 265–284. Springer (2006)
16. Flajolet, P., Martin, G.N.: Probabilistic counting algorithms for data base applications. *Journal of computer and system sciences* **31**(2), 182–209 (1985)
17. Gribonval, R., Blanchard, G., Keriven, N., Traonmilin, Y.: Compressive statistical learning with random feature moments. *Mathematical Statistics and Learning* **3**(2), 113–164 (2021)

18. Harder, F., Adamczewski, K., Park, M.: DP-MERF: Differentially private mean embeddings with randomFeatures for practical privacy-preserving data generation. In: Banerjee, A., Fukumizu, K. (eds.) Proceedings of the 24th International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research, vol. 130, pp. 1819–1827. PMLR (Apr 2021)
19. Huang, G., Huang, G.B., Song, S., You, K.: Trends in extreme learning machines: A review. *Neural Networks* **61**, 32–48 (Jan 2015). <https://doi.org/10.1016/j.neunet.2014.10.001>
20. Huang, G.B., Zhu, Q.Y., Siew, C.K.: Extreme learning machine: theory and applications. *Neurocomputing* **70**(1-3), 489–501 (2006)
21. Kenthapadi, K., Tran, T.T.: Pripearl: A framework for privacy-preserving analytics and reporting at linkedin. In: Proceedings of the 27th ACM International Conference on Information and Knowledge Management. pp. 2183–2191 (2018)
22. Keriven, N., Bourrier, A., Gribonval, R., Pérez, P.: Sketching for large-scale learning of mixture models. *Information and Inference: A Journal of the IMA* **7**(3), 447–508 (2018)
23. Keriven, N., Tremblay, N., Traonmilin, Y., Gribonval, R.: Compressive k-means. In: ICASSP (2017), <https://hal.inria.fr/hal-01386077/document>
24. Li, H., Xiong, L., Jiang, X.: Differentially private synthesization of multi-dimensional data using copula functions. In: Advances in database technology: proceedings. International conference on extending database technology. vol. 2014, p. 475. NIH Public Access (2014)
25. Liu, F., Huang, X., Chen, Y., Suykens, J.A.: Random features for kernel approximation: A survey on algorithms, theory, and beyond. *IEEE Transactions on Pattern Analysis & Machine Intelligence* (01), 1–1 (2021)
26. Misra, J., Gries, D.: Finding repeated elements. *Science of computer programming* **2**(2), 143–152 (1982)
27. de Montjoye, Y.A., Hidalgo, C.A., Verleysen, M., Blondel, V.D.: Unique in the crowd: The privacy bounds of human mobility. *Scientific reports* **3**, 1376 (2013)
28. de Montjoye, Y.A., Radaelli, L., Singh, V.K., et al.: Unique in the shopping mall: On the reidentifiability of credit card metadata. *Science* **347**(6221), 536–539 (2015)
29. Park, M., Vinaroz, M., Charusaie, M.A., Harder, F.: Polynomial magic! Hermite polynomials for private data generation. arXiv:2106.05042 [cs, stat] (Jun 2021)
30. Parzen, E.: On estimation of a probability density function and mode. *The annals of mathematical statistics* **33**(3), 1065–1076 (1962)
31. Qardaji, W., Yang, W., Li, N.: Prview: practical differentially private release of marginal contingency tables. In: Proceedings of the 2014 ACM SIGMOD international conference on Management of data. pp. 1435–1446 (2014)
32. Rahimi, A., Recht, B.: Random features for large-scale kernel machines. In: Advances in neural information processing systems. pp. 1177–1184 (2008)
33. Rahimi, A., Recht, B.: Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In: Advances in neural information processing systems. pp. 1313–1320 (2009)
34. Rudin, W.: *Fourier Analysis on Groups*. Interscience Publishers (1962)
35. Schellekens, V., Chatalic, A., Houssiau, F., de Montjoye, Y.A., Jacques, L., Gribonval, R.: Differentially private compressive k-means. In: ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 7933–7937. IEEE (2019)
36. Xie, L., Lin, K., Wang, S., Wang, F., Zhou, J.: Differentially private generative adversarial network. arXiv preprint arXiv:1802.06739 (2018)

37. Zhang, J., Cormode, G., Procopiuc, C.M., Srivastava, D., Xiao, X.: Privbayes: Private data release via bayesian networks. *ACM Transactions on Database Systems (TODS)* **42**(4), 1–41 (2017)
38. Zhang, R., Lan, Y., Huang, G.B., Xu, Z.B.: Universal Approximation of Extreme Learning Machine With Adaptive Growth of Hidden Nodes. *IEEE Transactions on Neural Networks and Learning Systems* **23**(2), 365–371 (Feb 2012). <https://doi.org/10.1109/TNNLS.2011.2178124>

A Proof of Theorem 1

Let J_Σ , the left-hand side of the inequality, the mean squared error between the empirical mean \bar{f} and the estimation from the sketch \tilde{f} . Denoting $X = (X_1, \dots, X_n)$, we have

$$\begin{aligned}
 J_\Sigma &= \mathbb{E}_{X, \xi} \left[\left(\frac{1}{n} \sum_{i=1}^n f(X_i) - \langle a, \frac{1}{n} (\sum_{i=1}^n \Phi(X_i) + \xi) \rangle \right)^2 \right] \\
 &= \mathbb{E}_{X, \xi} \left[\left(\frac{1}{n} \sum_{i=1}^n (f(X_i) - \langle a, \Phi(X_i) \rangle) - \frac{1}{n} \langle a, \xi \rangle \right)^2 \right] \\
 &\stackrel{(i)}{=} \mathbb{E}_X \left[\left(\frac{1}{n} \sum_{i=1}^n (f(X_i) - \langle a, \Phi(X_i) \rangle) \right)^2 \right] + \frac{1}{n^2} \mathbb{E}_\xi [\langle a, \xi \rangle^2] \\
 &\stackrel{(ii)}{=} \frac{n(n-1)}{n^2} \cdot (\mathbb{E}_X [f(X)] - \langle a, \mathbb{E}_X [\Phi(X)] \rangle)^2 \\
 &\quad + \frac{n}{n^2} \cdot \mathbb{E}_X [(f(X) - \langle a, \Phi(X) \rangle)^2] + \|a\|_2^2 \frac{\mathbb{V}[\xi]}{n^2}
 \end{aligned}$$

where we used in (i) the independence from ξ and X and the fact that $\mathbb{E}[\xi] = 0$, and in (ii) the fact that samples $(X_i)_{1 \leq i \leq n}$ are independent (and $\mathbb{V}[\cdot]$ denotes the variance of a random variable). Finally, we use Jensen's inequality (since $x \mapsto x^2$ is convex) to show that $(\mathbb{E}_X [f(X)] - \langle a, \mathbb{E}_X [\Phi(X)] \rangle)^2 \leq \mathbb{E}_X [(f(X) - \langle a, \Phi(X) \rangle)^2]$, which concludes the proof.

B M²M learning procedure

Input: Target function f , private data sketch $(\hat{z}_D, n + \zeta)$ with associated feature map Φ and noise level σ_ξ^2 , a priori distribution ψ , number of synthetic samples n_s , (optional) additional regularization $R = 1$.

Output: \hat{f} , an estimate for $\frac{1}{n} \sum_{i=1}^n f(x_i)$.

- 1 Get n_s synthetic training samples $\tilde{x}_i \sim_{i.i.d.} \psi$;
- 2 Set regularization parameter $\lambda = \sigma_\xi^2 / (|D| + \zeta) \cdot R$;
- 3 Get coefficients $a = \arg \min_\alpha J(\alpha)$, using the samples \tilde{x}_i as estimation for ψ ;
- 4 **return** $\hat{f} = \langle a, \hat{z}_D \rangle \approx \frac{1}{n} \sum_{i=1}^n f(x_i)$.

Algorithm 1: M²M learning procedure: Given a dataset sketch \hat{z}_D and a target function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the procedure estimates the empirical mean of f over D .