

Private Hypothesis Selection

Mark Bun^{*} Gautam Kamath[†] Thomas Steinke[‡] Zhiwei Steven Wu[§]

October 27, 2019

Abstract

We provide a differentially private algorithm for hypothesis selection. Given samples from an unknown probability distribution P and a set of m probability distributions \mathcal{H} , the goal is to output, in a ε -differentially private manner, a distribution from \mathcal{H} whose total variation distance to P is comparable to that of the best such distribution (which we denote by α). The sample complexity of our basic algorithm is $O\left(\frac{\log m}{\alpha^2} + \frac{\log m}{\alpha\varepsilon}\right)$, representing a minimal cost for privacy when compared to the non-private algorithm. We also can handle infinite hypothesis classes \mathcal{H} by relaxing to (ε, δ) -differential privacy.

We apply our hypothesis selection algorithm to give learning algorithms for a number of natural distribution classes, including Gaussians, product distributions, sums of independent random variables, piecewise polynomials, and mixture classes. Our hypothesis selection procedure allows us to generically convert a cover for a class to a learning algorithm, complementing known learning lower bounds which are in terms of the size of the packing number of the class. As the covering and packing numbers are often closely related, for constant α , our algorithms achieve the optimal sample complexity for many classes of interest. Finally, we describe an application to private distribution-free PAC learning.

1 Introduction

We consider the problem of *hypothesis selection*: given samples from an unknown probability distribution, select a distribution from some fixed set of candidates which is “close” to the unknown distribution in some appropriate distance measure. Such situations can arise naturally in a number of settings. For instance, we may have a number of different methods which work under various circumstances, which are not known in advance. One option is to run all the methods to generate a set of hypotheses, and pick the best from this set afterwards. Relatedly, an algorithm may branch its behavior based on a number of “guesses,” which will similarly result in a set of candidates, corresponding to the output at the end of each branch. Finally, if we know that the underlying distribution belongs to some (parametric) class, it is possible to essentially enumerate the class

^{*}Simons Institute for the Theory of Computing and Boston University. mbun@bu.edu. Supported by a Google Research Fellowship, as part of the Simons-Berkeley Research Fellowship program.

[†]Simons Institute for the Theory of Computing and University of Waterloo. g@csail.mit.edu. Supported as a Microsoft Research Fellow, as part of the Simons-Berkeley Research Fellowship program. Part of this work was completed while visiting Microsoft Research, Redmond.

[‡]IBM Research. phs@thomas-steinke.net. Part of this work completed while visiting the Simons Institute for the Theory of Computing at UC Berkeley.

[§]University of Minnesota, Twin Cities. zsw@umn.edu. Part of this work completed while visiting the Simons Institute for the Theory of Computing at UC Berkeley. Supported in part by a Google Faculty Research Award, a J.P. Morgan Faculty Award, and a Facebook Research Award.

(also known as a *cover*) to create a collection of hypotheses. Observe that this last example is quite general, and this approach can give generic learning algorithms for many settings of interest.

This problem of hypothesis selection has been extensively studied (see, e.g., [Yat85, DL96, DL97, DL01]), resulting in algorithms with a sample complexity which is *logarithmic* in the number of hypotheses. Such a mild dependence is critical, as it facilitates sample-efficient algorithms even when the number of candidates may be large. These initial works have triggered a great deal of study into hypothesis selection with additional considerations, including computational efficiency, understanding the optimal approximation factor, adversarial robustness, and weakening access to the hypotheses (e.g., [MS08, DDS12b, DK14, SOAJ14, AJOS14, DKK⁺16, AFJ⁺18, BKM19]).

However, in modern settings of data analysis, data may contain sensitive information about individuals. Some examples of such data include medical records, GPS location data, or private message transcripts. As such, we would like to perform statistical inference in these settings without revealing significant information about any particular individual's data. To this end, there have been many proposed notions of data privacy, but perhaps the gold standard is that of *differential privacy* [DMNS06]. Informally, differential privacy requires that, if a single datapoint in the dataset is changed, then the distribution over outputs produced by the algorithm should be similar (see Definition 2.4). Differential privacy has seen widespread adoption, including deployment by Apple [Dif17], Google [EPK14], and the US Census Bureau [DLS⁺17].

This naturally raises the question of whether one can perform hypothesis selection under the constraint of differential privacy, while maintaining a logarithmic dependence on the size of the cover. Such a tool would allow us to generically obtain private learning results for a wide variety of settings.

1.1 Results

Our main results answer this in the affirmative: we provide differentially private algorithms for selecting a good hypothesis from a set of distributions. The output distribution is competitive with the best distribution, and the sample complexity is bounded by the logarithm of the size of the set. The following is a basic version of our main result.

Theorem 1.1. *Let $\mathcal{H} = \{H_1, \dots, H_m\}$ be a set of probability distributions. Let $D = \{X_1, \dots, X_n\}$ be a set of samples drawn independently from an unknown probability distribution P . There exists an ε -differentially private algorithm (with respect to the dataset D) which has following guarantees. Suppose there exists a distribution $H^* \in \mathcal{H}$ such that $d_{TV}(P, H^*) \leq \alpha$. If $n = \Omega\left(\frac{\log m}{\alpha^2} + \frac{\log m}{\alpha\varepsilon}\right)$, then the algorithm will output a distribution $\hat{H} \in \mathcal{H}$ such that $d_{TV}(P, \hat{H}) \leq (3 + \zeta)\alpha$ with probability at least $9/10$, for any constant $\zeta > 0$. The running time of the algorithm is $O(nm^2)$.*

The sample complexity of this problem without privacy constraints is $O\left(\frac{\log m}{\alpha^2}\right)$, and thus the additional cost for ε -differential privacy is an additive $O\left(\frac{\log m}{\alpha\varepsilon}\right)$. We consider this cost to be minimal; in particular, the dependence on m is unchanged. Note that the running time of our algorithm is $O(nm^2)$ – we conjecture it may be possible to reduce this to $\tilde{O}(nm)$ as has been done in the non-private setting [DK14, SOAJ14, AJOS14, AFJ⁺18], though we have not attempted to perform this optimization. Regardless, our main focus is on the sample complexity rather than the running time, since any method for generic hypothesis selection requires $\Omega(m)$ time, thus precluding efficient algorithms when m is large. Note that the approximation factor of $(3 + \zeta)\alpha$ is effectively tight [DL01, MS08, BKM19]. Theorem 1.1 requires prior knowledge of the value of α , though we can use this to obtain an algorithm with similar guarantees which does not (Theorem 3.4).

It is possible to improve the guarantees of this algorithm in two ways (Theorem 4.1). First, if the distributions are nicely structured, the former term in the sample complexity can be reduced from $O(\log m/\alpha^2)$ to $O(d/\alpha^2)$, where d is a VC-dimension-based measure of the complexity of the collection of distributions. Second, if there are few hypotheses which are close to the true distribution, then we can pay only logarithmically in this number, as opposed to the total number of hypotheses. These modifications allow us to handle instances where m may be very large (or even infinite), albeit at the cost of weakening to approximate differential privacy to perform the second refinement. A technical discussion of our methods is in Section 1.2, our basic approach is covered in Section 3, and the version with all the bells and whistles appears in Section 4.

From Theorem 1.1, we immediately obtain Corollary 1.2 which applies when \mathcal{H} itself may not be finite, but admits a finite cover with respect to total variation distance.

Corollary 1.2. *Suppose there exists an α -cover \mathcal{C}_α of a set of distributions \mathcal{H} , and that we are given a set of samples $X_1, \dots, X_n \sim P$, where $d_{\text{TV}}(P, \mathcal{H}) \leq \alpha$. For any constant $\zeta > 0$, there exists an ε -differentially private algorithm (with respect to the input $\{X_1, \dots, X_n\}$) which outputs a distribution $H^* \in \mathcal{C}_\alpha$ such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, as long as*

$$n = \Omega \left(\frac{\log |\mathcal{C}_\alpha|}{\alpha^2} + \frac{\log |\mathcal{C}_\alpha|}{\alpha \varepsilon} \right).$$

Informally, this says that if a hypothesis class has an α -cover \mathcal{C}_α , then there is a private learning algorithm for the class which requires $O(\log |\mathcal{C}_\alpha|)$ samples. Note that our algorithm works even if the unknown distribution is only *close* to the hypothesis class. This is useful when we may have model misspecification, or when we require adversarial robustness. (We also give an extension of this algorithm which gives guarantees in the *semi-agnostic* learning model; see Section 3.3 for details.) The requirements for this theorem to apply are minimal, and thus it generically provides learning algorithms for a wide variety of hypothesis classes. That said, in non-private settings, the sample complexity given by this method is rather lossy: as an extreme example, there is no finite-size cover of univariate Gaussian distributions with unbounded parameters, so this approach does not give a finite-sample algorithm. That said, it is well-known that $O(1/\alpha^2)$ samples suffice to estimate a Gaussian in total variation distance. In the private setting, our theorem incurs a cost which is somewhat necessary: in particular, it is folklore that any pure ε -differentially private learning algorithm must pay a cost which is logarithmic in the packing number of the class (for completeness, see Lemma 5.1). Due to the relationship between packing and covering numbers (Lemma 5.2), this implies that up to a constant factor relaxation in the learning accuracy, our results are tight (Theorem 5.3). Further discussion appears in Sections 5.

Given Corollary 1.2, in Section 6, we derive new learning results for a number of classes. Our main applications are for d -dimensional Gaussian and product distributions. Informally, we obtain $\tilde{O}(d)$ sample algorithms for learning a product distribution and a Gaussian with known covariance (Corollaries 6.3 and 6.10), and an $\tilde{O}(d^2)$ algorithm for learning a Gaussian with unknown covariance (Corollary 6.11). These improve on recent results by Kamath, Li, Singhal, and Ullman [KLSU19] in two different ways. First, as mentioned before, our results are semi-agnostic, so we can handle when the distribution is only *close* to a product or Gaussian distribution. Second, our results hold for pure $(\varepsilon, 0)$ -differential privacy, which is a stronger notion than ε^2 -zCDP as considered in [KLSU19]. In this weaker model, they also obtained $\tilde{O}(d)$ and $\tilde{O}(d^2)$ sample algorithms, but the natural modifications to achieve ε -DP incur extra $\text{poly}(d)$ factors.¹ [KLSU19] also showed $\tilde{\Omega}(d)$

¹Roughly, this is due to the fact that the Laplace and Gaussian mechanism are based on ℓ_1 and ℓ_2 sensitivity, respectively, and that there is a \sqrt{d} -factor relationship between these two norms, in the worst case.

lower bounds for Gaussian and product distribution estimation in the even weaker model of (ε, δ) -differential privacy. Thus, our results show that the dimension dependence for these problems is unchanged for essentially any notion of differential privacy. In particular, our results show a previously-unknown separation between mean estimation of product distributions and non-product distributions under pure $(\varepsilon, 0)$ -differential privacy; see Remark 6.4.

We also apply Theorem 4.1 to obtain algorithms for learning Gaussians under (ε, δ) -differential privacy, with no bounds on the mean and variance parameters. More specifically, we provide algorithms for learning multivariate Gaussians with unknown mean and known covariance (Corollary 6.13), and univariate Gaussians with both unknown mean and variance (Corollary 6.15). For the former problem, we manage to avoid dependences which arise due to the application of advanced composition (similar to Remark 6.4).

To demonstrate the flexibility of our approach, we also give private learning algorithms for sums of independent random variables (Corollaries 6.20 and 6.22) and piecewise polynomials (Corollary 6.29). To the best of our knowledge, the former class of distributions has not been considered in the private setting, and we rely on covering theorems from the non-private literature. Private learning algorithms for the latter class, piecewise polynomials, have been studied by Diakonikolas, Hardt, and Schmidt [DHS15]. They provide sample and time efficient algorithms for histogram distributions (i.e., piecewise constant distributions), and claim similar results for general piecewise polynomials. Their method depends heavily on rather sophisticated algorithms for the non-private version of this problem [ADLS17]. In contrast, we can obtain comparable sample complexity bounds from just the existence of a cover and elementary VC dimension arguments, which we derive in a fairly self-contained manner.

We additionally give algorithms for learning mixtures of any coverable class (Corollary 6.32). In particular, this immediately implies algorithms for learning mixtures of Gaussians, product distributions, and all other classes mentioned above.

To conclude our applications, we discuss a connection to PAC learning (Corollary 6.34). It is known that the sample complexity of differentially private distribution-free PAC learning can be higher than that of non-private learning. However, this gap does not exist for distribution-specific learning, where the learning algorithm knows the distribution of (unlabeled) examples, as both sample complexities are characterized by VC dimension. Private hypothesis selection allows us to address an intermediate situation where the distribution of unlabeled examples is not known exactly, but is known to come (approximately) from a class of distributions. When this class has a small cover, we are able to recover sample complexity guarantees for private PAC learning which are comparable to the non-private case.

1.2 Techniques

Non-privately, most algorithms for hypothesis selection involve a tournament-style approach. We conduct a number of pairwise comparisons between distributions, which may either have a winner and a loser, or may be declared a draw. Intuitively, a distribution will be declared the winner of a comparison if it is much closer than the alternative to the unknown distribution, and a tie will be declared if the two distributions are comparably close. The algorithm will output any distribution which never loses a comparison. A single comparison between a pair of hypotheses requires $O(1/\alpha^2)$ samples, and a Chernoff plus union bound argument over the $O(m^2)$ possible comparisons increases the sample complexity to $O(\log m/\alpha^2)$. In fact, we can use uniform convergence arguments to reduce this sample complexity to $O(d/\alpha^2)$, where d is the VC dimension of the $2^{\binom{m}{2}}$ sets (the ‘‘Scheffé’’ sets) defined by the subsets of the domain where the PDF of one distribution dominates another. Crucially, we must reuse the same set of samples for all comparisons to avoid paying polynomially

in the number of hypotheses.

A private algorithm for this problem requires additional care. Since a single comparison is based on the number of samples which fall into a particular subset of the domain, the sensitivity of the underlying statistic is low, and thus privacy may seem easily achievable at first glance. However, the challenge comes from the fact that the same samples are reused for all pairwise comparisons, thus greatly increasing the sensitivity: changing a single datapoint could flip the result of every comparison! In order to avoid this pitfall, we instead carefully construct a score function for each hypothesis, namely, the minimum number of points that must be changed to cause the distribution to lose any comparison. For this to be a useful score function, we must show that the best hypothesis will win all of its comparisons by a large margin. We can then use the Exponential Mechanism [MT07] to select a distribution with high score.

Further improvements can be made if we are guaranteed that the number of “good” hypotheses (i.e., those that have total variation distance from the true distribution bounded by $(3 + \zeta)\alpha$) is at most some parameter k , and if we are willing to relax to approximate differential privacy. The parameter k here is related to the doubling dimension of the hypothesis class with respect to total variation distance. If we randomly assign the hypotheses to $\Omega(k^2)$ buckets, with high probability, no bucket will contain more than one good hypothesis. We can identify a bucket containing a good hypothesis using a similar method based on the exponential mechanism as described above. Moreover, since we are likely to only have one “good” hypothesis in the chosen bucket, this implies a significant gap between the best and second-best scores in that bucket. This allows us to use stability-based techniques [DL09, TS13], and in particular the GAP-MAX algorithm of Bun, Dwork, Rothblum, and Steinke [BDRS18], to identify an accurate distribution.

1.3 Related Work

Our main result builds on a long line of work on non-private hypothesis selection. One starting point for the particular style of approach we consider here is [Yat85], which was expanded on in [DL96, DL97, DL01]. Since then, there has been study into hypothesis selection under additional considerations, including computational efficiency, understanding the optimal approximation factor, adversarial robustness, and weakening access to the hypotheses [MS08, DDS12b, DK14, SOAJ14, AJOS14, DKK⁺16, AFJ⁺18, BKM19]. Our private algorithm examines the same type of problem, with the additional constraint of differential privacy.

There has recently been a great deal of interest in differentially private distribution learning. In the central model, most relevant are [DHS15], which gives algorithms for learning structured univariate distributions, and [KV18, KLSU19], which focus on learning Gaussians and binary product distributions. [CWZ19] also studies private statistical parameter estimation. Privately learning mixtures of Gaussians was considered in [NRS07, KSSU19]. The latter paper (which is concurrent with the present work) gives a computationally efficient algorithm for the problem, but with a worse sample complexity, and incomparable accuracy guarantees (they require a separation condition, and perform clustering and parameter estimation, while we do proper learning). [BNSV15] give an algorithm for learning distributions in Kolmogorov distance. Upper and lower bounds for learning the mean of a product distribution over the hypercube in ℓ_∞ -distance include [BDMN05, BUV14, DMNS06, SU17]. [AKSZ18] focuses on estimating properties of a distribution, rather than the distribution itself. [Smi11] gives an algorithm which allows one to estimate asymptotically normal statistics with optimal convergence rates, but no finite sample complexity guarantees. There has also been a great deal of work on distribution learning in the local model of differential privacy [DJW13, WHW⁺16, KBR16, ASZ19, DR18, JKMW18, YB18, GRS19].

Non-privately, there has been a significant amount of work on learning specific classes of distri-

butions. The PAC-style formulation of the problem we consider originated in [KMR⁺94]. While learning Gaussians and product distributions can be considered folklore at this point, some of the other classes we learn have enjoyed more recent study. For instance, learning sums of independent random variables was recently considered in [DDS12b] toward the problem of learning Poisson Binomial Distributions (PBDs). Since then, there has been additional work on learning PBDs and various generalizations [DKT15, DDKT16, DKS16b, DKS16c, DKS16a, DLS18].

Piecewise polynomials are a highly-expressive class of distributions, and they can be used to approximate a number of other univariate distribution classes, including distributions which are multi-modal, concave, convex, log-concave, monotone hazard rate, Gaussian, Poisson, Binomial, and more. Algorithms for learning such classes are considered in a number of papers, including [DDS12a, CDSS14a, CDSS14b, ADK15, ADLS17].

There has also been a great deal of work on learning mixtures of distribution classes, particularly mixtures of Gaussians. There are many ways the objective of such a problem can be defined, including clustering [Das99, DS00, AK01, VW02, AM05, CR08b, CR08a, KK10, AS12, RV17, HL18, DKS18, KSS18], parameter estimation [KMV10, MV10, BS10, HK13, ABG⁺14, BCMV14, HP15, GHK15, XHM16, DTZ17, ABDH⁺18], proper learning [FOS06, FOS08, DK14, SOAJ14, DKK⁺16, LS17], and improper learning [CDSS14a]. Our work falls into the line on proper learning: the algorithm is given a set of samples from a mixture of Gaussians, and must output a mixture of Gaussians which is close in total variation distance.

1.4 Organization

We begin in Section 2 with preliminaries. In Section 3, we give a basic algorithm for private hypothesis selection, via the exponential mechanism. In Section 4, we extend this approach in two ways: by using VC dimension arguments to reduce the sample complexity for sets of hypotheses with additional structure, and combining this with a GAP-MAX algorithm to achieve non-trivial guarantees for infinite hypothesis classes. Section 5 shows that our approach leads to algorithms which essentially match lower bounds for most distribution classes (in the constant α regime). We consider applications in Section 6: through a combination of arguments about covers and VC dimension, we derive algorithms for learning a number of classes of distributions, as well as describe an application to private PAC learning. Finally, we conclude in Section 7 with open questions.

2 Preliminaries

We start with some preliminaries and definitions.

Definition 2.1. *The total variation distance or statistical distance between P and Q is defined as*

$$d_{\text{TV}}(P, Q) = \max_{S \subseteq \Omega} P(S) - Q(S) = \frac{1}{2} \int_{x \in \Omega} |P(x) - Q(x)| dx = \frac{1}{2} \|P - Q\|_1 \in [0, 1].$$

Moreover, if \mathcal{H} is a set of distributions over a common domain, we define $d_{\text{TV}}(P, \mathcal{H}) = \inf_{H \in \mathcal{H}} d_{\text{TV}}(P, H)$.

Throughout this paper, we consider packings and coverings of sets of distributions with respect to total variation distance.

Definition 2.2. *A γ -cover of a set of distributions \mathcal{H} is a set of distributions \mathcal{C}_γ , such that for every $H \in \mathcal{H}$, there exists some $P \in \mathcal{C}_\gamma$ such that $d_{\text{TV}}(P, H) \leq \gamma$.*

A γ -packing of a set of distributions \mathcal{H} is a set of distributions $\mathcal{P}_\gamma \subseteq \mathcal{H}$, such that for every pair of distributions $P, Q \in \mathcal{P}_\gamma$, we have that $d_{\text{TV}}(P, Q) \geq \gamma$.

In this paper, we present semi-agnostic learning algorithms.

Definition 2.3. *An algorithm is said to be an α -semi-agnostic learner for a class \mathcal{H} if it has the following guarantees. Suppose we are given $X_1, \dots, X_n \sim P$, where $d_{\text{TV}}(P, \mathcal{H}) \leq \text{OPT}$. The algorithm must output some distribution \hat{H} such that $d_{\text{TV}}(P, \hat{H}) \leq c \cdot \text{OPT} + O(\alpha)$, for some constant $c \geq 1$. If $c = 1$, then the algorithm is said to be agnostic.*

Now we define differential privacy. We say that D and D' are neighboring datasets, denoted $D \sim D'$, if D and D' differ by at most one observation. Informally, differential privacy requires that the algorithm has close output distributions when run on any pair of neighboring datasets. More formally:

Definition 2.4 ([DMNS06]). *A randomized algorithm $T : X^* \rightarrow \mathcal{R}$ is (ϵ, δ) -differentially private if for all $n \geq 1$, for all neighboring datasets $D, D' \in X^n$, and for all events $S \subseteq \mathcal{R}$,*

$$\Pr[T(D) \in S] \leq e^\epsilon \Pr[T(D') \in S] + \delta.$$

If $\delta = 0$, we say that T is ϵ -differentially private.

We will also use the related notion of concentrated differential privacy:

Definition 2.5 ([DR16, BS16]). *A randomized algorithm $T : X^* \rightarrow \mathcal{R}$ satisfies ρ -zero-concentrated differential privacy if for all $n \geq 1$, for all neighboring datasets $D, D' \in X^n$, and for all $\alpha \in (1, \infty)$,*

$$R_\alpha(M(D) || M(D')) \leq \rho \alpha,$$

where $R_\alpha(M(D) || M(D'))$ is the α -Rényi divergence between $M(D)$ and $M(D')$.²

The exponential mechanism [MT07] is a powerful ϵ -differentially private mechanism for selecting an approximately best outcome from a set of alternatives, where the quality of an outcome is measured by a score function relating each alternative to the underlying dataset. Letting \mathcal{R} be the set of possible outcomes, a score function $q : X^* \times \mathcal{R} \rightarrow \mathbb{R}$ maps each pair consisting of a dataset and an outcome to a real-valued score. The exponential mechanism \mathcal{M}_E instantiated with a dataset D , a score function q , and a privacy parameter ϵ selects an outcome r in \mathcal{R} with probability proportional to $\exp(\epsilon q(D, r) / (2\Delta(q)))$, where $\Delta(q)$ is the sensitivity of the score function defined as

$$\Delta(q) = \max_{r \in \mathcal{R}, D \sim D'} |q(D, r) - q(D', r)|.$$

Theorem 2.6 ([MT07]). *For any input dataset D , score function q and privacy parameter $\epsilon > 0$, the exponential mechanism $\mathcal{M}_E(D, q, \epsilon)$ is ϵ -differentially private, and with probability at least $1 - \beta$, selects an outcome $r \in \mathcal{R}$ such that*

$$q(D, r) \geq \max_{r' \in \mathcal{R}} q(D, r') - \frac{2\Delta(q) \log(|\mathcal{R}|/\beta)}{\epsilon}.$$

3 A First Method for Private Hypothesis Selection

In this section, we present our first algorithm for private hypothesis selection and obtain the following result.

²Given two probability distributions P, Q over Ω , $R_\alpha(P || Q) = \frac{1}{\alpha-1} \log(\sum_{x \in \Omega} P(x)^\alpha Q(x)^{1-\alpha})$.

Theorem 1.1. Let $\mathcal{H} = \{H_1, \dots, H_m\}$ be a set of probability distributions. Let $D = \{X_1, \dots, X_n\}$ be a set of samples drawn independently from an unknown probability distribution P . There exists an ε -differentially private algorithm (with respect to the dataset D) which has following guarantees. Suppose there exists a distribution $H^* \in \mathcal{H}$ such that $d_{\text{TV}}(P, H^*) \leq \alpha$. If $n = \Omega\left(\frac{\log m}{\alpha^2} + \frac{\log m}{\alpha\varepsilon}\right)$, then the algorithm will output a distribution $\hat{H} \in \mathcal{H}$ such that $d_{\text{TV}}(P, \hat{H}) \leq (3+\zeta)\alpha$ with probability at least $9/10$, for any constant $\zeta > 0$. The running time of the algorithm is $O(nm^2)$.

Note that the sample complexity bound above scales logarithmically with the size of the hypothesis class. In Section 4, we will provide a stronger result (which subsumes the present one as a special case) that can handle certain infinite hypothesis classes. For sake of exposition, we begin in this section with the basic algorithm.

3.1 Pairwise Comparisons

We first present a subroutine which compares two hypothesis distributions. Let H and H' be two distributions over domain \mathcal{X} and consider the following set, which is called the *Scheffé set*:

$$\mathcal{W}_1 = \{x \in \mathcal{X} \mid H(x) > H'(x)\}$$

Define $p_1 = H(\mathcal{W}_1)$, $p_2 = H'(\mathcal{W}_1)$, and $\tau = P(\mathcal{W}_1)$ to be the probability masses that H , H' , and P place on \mathcal{W}_1 , respectively. It follows that $p_1 > p_2$ and $p_1 - p_2 = d_{\text{TV}}(H, H')$.³

Algorithm 1: PAIRWISE CONTEST: $\text{PC}(H, H', D, \zeta, \alpha)$

Input: Two hypotheses H and H' , input dataset D of size n drawn i.i.d. from target distribution P , approximation parameter $\zeta > 0$, and accuracy parameter $\alpha \in (0, 1)$.
Initialize: Compute the fraction of points that fall into \mathcal{W}_1 : $\hat{\tau} = \frac{1}{n} |\{x \in D \mid x \in \mathcal{W}_1\}|$.
If $p_1 - p_2 \leq (2 + \zeta)\alpha$, return “Draw”.
Else If $\hat{\tau} > p_1 - (1 + \zeta/2)\alpha$, return H as the winner.
Else If $\hat{\tau} < p_2 + (1 + \zeta/2)\alpha$, return H' as the winner.
Else return “Draw”.

Now consider the following function of this ordered pair of hypotheses:

$$\Gamma_\zeta(H, H', D) = \begin{cases} n & \text{if } p_1 - p_2 \leq (2 + \zeta)\alpha; \\ n \cdot \max\{0, \hat{\tau} - (p_2 + (1 + \zeta/2)\alpha)\} & \text{otherwise.} \end{cases}$$

When the two hypotheses are sufficiently far apart (i.e., $d_{\text{TV}}(H, H') > (2 + \zeta)\alpha$), $\Gamma_\zeta(H, H', D)$ is essentially the number of points one needs to change in D to make H' the winner.

Lemma 3.1. Let P, H, H' be distributions as above. With probability at least $1 - 2\exp(-n\zeta^2\alpha^2/8)$ over the random draws of D from P^n , $\hat{\tau}$ satisfies $|\hat{\tau} - \tau| < \zeta\alpha/4$, and if $d_{\text{TV}}(P, H) \leq \alpha$, then $\Gamma_\zeta(H, H', D) > \zeta\alpha n/4$.

Proof. By applying Hoeffding’s inequality, we know that with probability at least $1 - 2\exp(-n\zeta^2\alpha^2/8)$, $|\tau - \hat{\tau}| < \zeta\alpha/4$. We condition on this event for the remainder of the proof. Consider the following two cases. In the first case, suppose that $p_1 - p_2 \leq (2 + \zeta)\alpha$. Then we know that $\Gamma_\zeta(H, H', D) = n > \alpha n$. In the second case, suppose that $p_1 - p_2 > (2 + \zeta)\alpha$. Since $d_{\text{TV}}(P, H) \leq \alpha$, we know that $|p_1 - \tau| \leq \alpha$, and so $|p_1 - \hat{\tau}| < (1 + \zeta/4)\alpha$. Since $p_1 > p_2 + (2 + \zeta)\alpha$, we also have $\hat{\tau} > p_2 + (1 + 3\zeta/4)\alpha$. It follows that $\Gamma_\zeta(H, H', D) = n(\hat{\tau} - (p_2 + (1 + \zeta/2)\alpha)) > \zeta\alpha n/4$. This completes the proof. \square

³For simplicity of our exposition, we will assume that we can evaluate the two quantities p_1 and p_2 exactly. In general, we can estimate these quantities to arbitrary accuracy, as long as, for each hypothesis H , we can evaluate the density of each point under H and also draw samples from H .

3.2 Selection via Exponential Mechanism

In light of the definition of the pairwise comparison defined above, we consider the following score function $S: \mathcal{H} \times \mathcal{X}^n$, such that for any $H_j \in \mathcal{H}$ and dataset D ,

$$S(H_j, D) = \min_{H_k \in \mathcal{H}} \Gamma_\zeta(H_j, H_k, D). \quad (1)$$

Roughly speaking, $S(H_j, D)$ is the minimum number of points required to change in D in order for H_j to lose at least one pairwise contest against a different hypothesis. When the hypothesis H_j is very close to every other distribution, such that all pairwise contests return “Draw,” then the score will be n .

Algorithm 2: PRIVATE HYPOTHESIS SELECTION: $\text{PHS}(\mathcal{H}, D, \varepsilon)$

Input: Dataset D , a collection of hypotheses $\mathcal{H} = \{H_1, \dots, H_m\}$, privacy parameter ε .

Output a random hypothesis $\hat{H} \in \mathcal{H}$ such that for each H_j

$$\Pr[\hat{H} = H_j] \propto \exp\left(\frac{S(H_j, D)}{2\varepsilon}\right)$$

where $S(H_j, D)$ is defined in (1).

Lemma 3.2 (Privacy). *For any $\varepsilon > 0$ and collection of hypotheses \mathcal{H} , the algorithm $\text{PHS}(\mathcal{H}, \cdot, \varepsilon)$ satisfies ε -differential privacy.*

Proof. First, observe that for any pairs of hypotheses H_j, H_k , $\Gamma_\zeta(H_j, H_k, \cdot)$ has sensitivity 1. As a result, the score function S is also 1-sensitive. Then the result directly follows from the privacy guarantee of the exponential mechanism (Theorem 2.6). \square

Lemma 3.3 (Utility). *Fix any $\alpha, \beta \in (0, 1)$, and $\zeta > 0$. Suppose that there exists $H^* \in \mathcal{H}$ such that $d_{\text{TV}}(P, H^*) \leq \alpha$. Then with probability $1 - \beta$ over the sample D and the algorithm PHS , we have that $\text{PHS}(\mathcal{H}, D)$ outputs an hypothesis \hat{H} such that $d_{\text{TV}}(P, \hat{H}) \leq (3 + \zeta)\alpha$, as long as the sample size satisfies*

$$n \geq \frac{8 \ln(4m/\beta)}{\zeta^2 \alpha^2} + \frac{8 \ln(2m/\beta)}{\zeta \alpha \varepsilon}.$$

Proof. First, consider the m pairwise contests between H^* and every candidate in \mathcal{H} . Let $\mathcal{W}_j = \{x \in \mathcal{X} \mid H_j(x) > H^*(x)\}$ be the collection of Scheffé sets. For any event $W \subseteq \mathcal{X}$, let $\hat{P}(W)$ denote the empirical probability of event W on the dataset D . By Lemma 3.1 and an application of the union bound, we know that with probability at least $1 - 2m \exp(-n\zeta^2\alpha^2/8)$ over the draws of D , $|P(\mathcal{W}_j) - \hat{P}(\mathcal{W}_j)| \leq \zeta\alpha/4$ and $\Gamma_\zeta(H^*, H_j, D) > \zeta\alpha n/4$ for all $H_j \in \mathcal{H}$. In particular, the latter event implies that $S(H^*, D) > \zeta\alpha n/4$.

Next, by the utility guarantee of the exponential mechanism (Theorem 2.6), we know that with probability at least $1 - \beta/2$, the output hypothesis satisfies

$$S(\hat{H}, D) \geq S(H^*, D) - \frac{2 \ln(2m/\beta)}{\varepsilon} > \zeta\alpha n/4 - \frac{2 \ln(2m/\beta)}{\varepsilon}.$$

Then as long as $n \geq \frac{8 \ln(4m/\beta)}{\zeta^2 \alpha^2} + \frac{8 \ln(2m/\beta)}{\zeta \alpha \varepsilon}$, we know that with probability at least $1 - \beta$, $S(\hat{H}, D) > 0$.

Let us condition on this event, which implies that $\Gamma_\zeta(\hat{H}, H^*, D) > 0$. We will now show that $d_{\text{TV}}(\hat{H}, H^*) \leq (2 + \zeta)\alpha$, which directly implies that $d_{\text{TV}}(\hat{H}, P) \leq (3 + \zeta)\alpha$ by the triangle inequality. Suppose to the contrary that $d_{\text{TV}}(\hat{H}, H^*) > (2 + \zeta)\alpha$. Then by the definition of Γ_ζ , $\hat{P}(\hat{\mathcal{W}}) > H^*(\hat{\mathcal{W}}) + (1 + \zeta/2)\alpha$, where $\hat{\mathcal{W}} = \{x \in \mathcal{X} \mid \hat{H}(x) > H^*(x)\}$. Since $|P(\hat{\mathcal{W}}) - \hat{P}(\hat{\mathcal{W}})| \leq \zeta\alpha/4$, we have $P(\hat{\mathcal{W}}) > H^*(\hat{\mathcal{W}}) + (1 + \zeta/4)\alpha$, which is a contradiction to the assumption that $d_{\text{TV}}(P, H^*) \leq \alpha$. \square

3.3 Obtaining a Semi-Agnostic Algorithm

Theorem 1.1 shows that given a hypothesis class \mathcal{H} and samples from an unknown distribution P , we can privately find a distribution $\hat{H} \in \mathcal{H}$ with $d_{\text{TV}}(P, \hat{H}) \leq (3 + \zeta)\alpha$ *provided* that we know $d_{\text{TV}}(P, \mathcal{H}) \leq \alpha$. But what if we are not promised that P is itself close to \mathcal{H} ? We would like to design a private hypothesis selection algorithm for the more general semi-agnostic setting, where for any value of $\text{OPT} := d_{\text{TV}}(P, \mathcal{H})$, we are able to privately identify a distribution $\hat{H} \in \mathcal{H}$ with $d_{\text{TV}}(P, \hat{H}) \leq c \cdot \text{OPT} + \alpha$ for some universal constant c . Our goal will be to do this with sample complexity which is still logarithmic in $|\mathcal{H}|$.

Our strategy for handling this more general setting is by a reduction to that of Theorem 1.1. We run that algorithm $T = O(\log(1/\alpha))$ times, doubling the choice of α in each run and producing a sequence of candidate hypotheses H_1, \dots, H_T . By the guarantees of Theorem 1.1, there is some candidate H_t with $d_{\text{TV}}(P, H_t) \leq 2(3 + \zeta)\text{OPT}$. The remaining task is to approximately select the best candidate from H_1, \dots, H_T . This is done by implementing a private version of the Scheffé tournament which is itself semi-agnostic, but has a very poor (quadratic) dependence on the number of candidates T .

We prove the following result, which gives a semi-agnostic learner whose sample complexity is comparable to that of Theorem 1.1.

Theorem 3.4. *Let $\alpha, \beta, \varepsilon \in (0, 1)$, and $\zeta > 0$ be a constant. Let \mathcal{H} be a set of m distributions and let P be a distribution with $d_{\text{TV}}(P, \mathcal{H}) = \text{OPT}$. There is an ε -differentially private algorithm which takes as input n samples from P and with probability at least $1 - \beta$, outputs a distribution $\hat{H} \in \mathcal{H}$ with $d_{\text{TV}}(P, \hat{H}) \leq 18(3 + \zeta)\text{OPT} + \alpha$, as long as*

$$n \geq O\left(\frac{\log(m/\beta) + \log \log(1/\alpha)}{\alpha^2} + \frac{\log m + \log^2(1/\alpha) \cdot (\log(1/\beta) + \log \log(1/\alpha))}{\alpha \varepsilon}\right).$$

As discussed above, the algorithm relies on the following variant with a much worse dependence on m .

Lemma 3.5. *Let $\alpha, \beta, \varepsilon \in (0, 1)$. There is an ε -differentially private algorithm which takes as input n samples from P and with probability at least $1 - \beta$, outputs a distribution $\hat{H} \in \mathcal{H}$ with $d_{\text{TV}}(P, \hat{H}) \leq 9\text{OPT} + \alpha$, as long as*

$$n \geq O\left(\frac{\log(m/\beta)}{\alpha^2} + \frac{m^2 \log(m/\beta)}{\alpha \varepsilon}\right).$$

Proof sketch. We use a different variation of the Scheffé tournament which appears in [DL01]. Non-privately, the algorithm works as follows. For every pair of hypotheses $H, H' \in \mathcal{H}$ with Scheffé set $\mathcal{W}_{H,H'} = \{x \in \mathcal{X} \mid H(x) > H'(x)\}$, let $H(\mathcal{W}_{H,H'})$, $H'(\mathcal{W}_{H,H'})$, and $P(\mathcal{W}_{H,H'})$ denote the probability masses of H, H', P on $\mathcal{W}_{H,H'}$, respectively. Moreover, let $\hat{P}(\mathcal{W}_{H,H'})$ denote the fraction of points in the input sample D which lie in $\mathcal{W}_{H,H'}$. We declare H to be the winner of the pairwise contest between H and H' if $|H(\mathcal{W}_{H,H'}) - \hat{P}(\mathcal{W}_{H,H'})| < |H'(\mathcal{W}_{H,H'}) - \hat{P}(\mathcal{W}_{H,H'})|$. Otherwise, we declare H' to be the winner. The algorithm outputs the hypothesis \hat{H} which wins the most pairwise contests (breaking ties arbitrarily).

To make this algorithm ε -differentially private, we replace $\hat{P}(\mathcal{W}_{H,H'})$ in each pairwise contest with the $(\varepsilon/\binom{m}{2})$ -differentially private estimate $c_{H,H'} = \hat{P}(\mathcal{W}_{H,H'}) + \text{Lap}(\binom{m}{2}/\varepsilon n)$. By the composition guarantees of differential privacy, the algorithm as a whole is ε -differentially private.

The analysis of Devroye and Lugosi [DL01, Theorem 6.2] shows that the (private) Scheffé tournament outputs a hypothesis \hat{H} with

$$d_{\text{TV}}(\hat{H}, P) \leq 9\text{OPT} + 16 \max_{H, H' \in \mathcal{H}} |P(\mathcal{W}_{H,H'}) - c_{H,H'}|.$$

Fix an arbitrary pair H, H' . A Chernoff bound shows that $|P(\mathcal{W}_{H,H'}) - \hat{P}(\mathcal{W}_{H,H'})| \leq \alpha/32$ with probability at least $1 - \beta/(2m^2)$ as long as $n \geq O(\ln(m/\beta)/\alpha^2)$. Moreover, properties of the Laplace distribution guarantee $|c_{H,H'} - \hat{P}(\mathcal{W}_{H,H'})| \leq \alpha/32$ with probability at least $1 - \beta/(2m^2)$ as long as $n \geq O(m^2 \log(m/\beta)/\alpha\epsilon)$. The triangle inequality and a union bound over all pairs H, H' complete the proof. \square

Proof of Theorem 3.4. We now combine the private hypothesis selection algorithm of Theorem 1.1 with the expensive semi-agnostic learner of Lemma 3.5 to prove Theorem 3.4. Define sequences $\alpha_1 = \alpha/126, \alpha_2 = 2\alpha/126, \dots, \alpha_T = 2^{T-1}\alpha/126$ and $\epsilon_1 = \epsilon/4, \epsilon_2 = \epsilon/8, \dots, \epsilon_T = 2^{-(T+1)}\epsilon$ for $T = \lceil \log_2(1/\alpha) \rceil + 1$. For each $t = 1, \dots, T$, let H_t denote the outcome of a run of Algorithm 2 using accuracy parameter α_t and privacy parameter ϵ_t . Finally, use the algorithm of Lemma 3.5 to select a hypothesis from H_0, \dots, H_T using accuracy parameter α and privacy parameter $\epsilon/2$.

Privacy of this algorithm follows immediately from composition of differential privacy. We now analyze its sample complexity guarantee. By Lemma 3.3, we have that all T runs of Algorithm 2 succeed simultaneously with probability at least $1 - \beta/2$ as long as

$$n \geq O\left(\frac{\log(m/\beta) + \log \log(1/\alpha)}{\alpha^2} + \frac{\log(m/\beta) + \log \log(1/\alpha)}{\alpha\epsilon}\right).$$

Condition on this event occurring. Recall that success of run t of Algorithm 2 means that if $\text{OPT} \in (\alpha_{t-1}, \alpha_t]$, then $d_{\text{TV}}(P, H_t) \leq (3 + \zeta)\alpha_t \leq 2(3 + \zeta)\text{OPT}$. Meanwhile, if $\text{OPT} \leq \alpha_1 = \alpha/126$, then we have $d_{\text{TV}}(P, H_1) \leq \alpha/18$. Hence, regardless of the value of OPT , there exists a run t such that $d_{\text{TV}}(P, H_t) \leq 2(3 + \zeta)\text{OPT} + \alpha/18$. The algorithm of Lemma 3.5 is now, with probability at least $1 - \beta/2$, able to select a hypothesis \hat{H} with $d_{\text{TV}}(P, \hat{H}) \leq 9d_{\text{TV}}(P, H_t) + \alpha/2 \leq 18(3 + \zeta)\text{OPT} + \alpha$ as long as

$$n \geq O\left(\frac{\log(1/\beta) + \log \log(1/\alpha)}{\alpha^2} + \frac{\log^2(1/\alpha) \cdot (\log(1/\beta) + \log \log(1/\alpha))}{\alpha\epsilon}\right).$$

This gives the asserted sample complexity guarantee. \square

4 An Advanced Method for Private Hypothesis Selection

In Section 3, we provided a simple algorithm whose sample complexity grows logarithmically in the size of the hypothesis class. We now demonstrate that this dependence can be improved and, indeed, we can handle infinite hypothesis classes given that their VC dimension is finite and that the cover has small doubling dimension.

To obtain this improved dependence on the hypothesis class size, we must make two improvements to the analysis and algorithm. First, rather than applying a union bound over all the pairwise contests to analyse the tournament, we use a uniform convergence bound in terms of the VC dimension of the Scheffé sets. Second, rather than use the exponential mechanism to select a hypothesis, we use a “GAP-MAX” algorithm [BDRS18]. This takes advantage of the fact that, in many cases, even for infinite hypothesis classes, only a handful of hypotheses will have high scores. The GAP-MAX algorithm need only pay for the hypotheses that are close to optimal. To exploit this, we must move to a relaxation of pure differential privacy which is not subject to strong packing lower bounds (as we describe in Section 5). Specifically, we consider approximate differential privacy, although results with an improved dependence are also possible under various variants of concentrated differential privacy [DR16, BS16, Mir17, BDRS18].

Theorem 4.1. *Let \mathcal{H} be a set of probability distributions on \mathcal{X} . Let d be the VC dimension of the set of functions $f_{H,H'} : \mathcal{X} \rightarrow \{0,1\}$ defined by $f_{H,H'}(x) = 1 \iff H(x) > H'(x)$ where $H, H' \in \mathcal{H}$. There exists a (ε, δ) -differentially private algorithm which has following guarantee. Let $D = \{X_1, \dots, X_n\}$ be a set of private samples drawn independently from an unknown probability distribution P . Let $k = |\{H \in \mathcal{H} : d_{\text{TV}}(H, P) \leq 7\alpha\}|$. Suppose there exists a distribution $H^* \in \mathcal{H}$ such that $d_{\text{TV}}(P, H^*) \leq \alpha$. If $n = \Omega\left(\frac{d + \log(1/\beta)}{\alpha^2} + \frac{\log(k/\beta) + \min\{\log|\mathcal{H}|, \log(1/\delta)\}}{\alpha\varepsilon}\right)$, then the algorithm will output a distribution $\hat{H} \in \mathcal{H}$ such that $d_{\text{TV}}(P, \hat{H}) \leq 7\alpha$ with probability at least $1 - \beta$.*

Alternatively, we can demand that the algorithm be $\frac{1}{2}\varepsilon^2$ -concentrated differentially private if $n = \Omega\left(\frac{d + \log(1/\beta)}{\alpha^2} + \frac{\log(k/\beta) + \sqrt{\log|\mathcal{H}|}}{\alpha\varepsilon}\right)$.

Comparing Theorem 4.1 to Theorem 1.1, we see that the first (non-private) $\log|\mathcal{H}|$ term is replaced by the VC dimension d and the second (private) $\log|\mathcal{H}|$ term is replaced by $\log k + \log(1/\delta)$. Here k is a measure of the “local” size of the hypothesis class \mathcal{H} ; its definition is similar to that of the doubling dimension of the hypothesis class under total variation distance.

We note that the $\log(1/\delta)$ term could be large, as the privacy failure probability δ should be cryptographically small. Thus our result includes statements for pure differential privacy (by using the other term in the minimum with $\delta = 0$) and also concentrated differential privacy. Note that, since d and $\log k$ can be upper-bounded by $O(\log|\mathcal{H}|)$, this result supercedes the guarantees of Theorem 1.1.

4.1 VC Dimension

We begin by reviewing the definition of Vapnik-Chervonenkis (VC) dimension and its properties.

Definition 4.2 (VC dimension [VC74]). *Let \mathcal{F} be a set of functions $f : \mathcal{X} \rightarrow \{0,1\}$. The VC dimension of \mathcal{F} is defined to be the largest d such that there exist $x_1, \dots, x_d \in \mathcal{X}$ and $f_1, \dots, f_{2^d} \in \mathcal{F}$ such that for all $1 \leq i < j \leq 2^d$ there exists $1 \leq k \leq d$ such that $f_i(x_k) \neq f_j(x_k)$.*

For our setting, we must extend the definition of VC dimension from function families to hypothesis classes.

Definition 4.3 (VC dimension of hypothesis class). *Let \mathcal{H} be a set of probability distributions on a space \mathcal{X} . For $H, H' \in \mathcal{H}$, define $f_{H,H'} : \mathcal{X} \rightarrow \{0,1\}$ by $f_{H,H'}(x) = 1 \iff H(x) > H'(x)$. Define $\mathcal{F}(\mathcal{H}) = \{f_{H,H'} : H, H' \in \mathcal{H}\}$. We define the VC dimension of \mathcal{H} to be the VC dimension of $\mathcal{F}(\mathcal{H})$.⁴*

The key property of VC dimension is the following uniform convergence bound, which we use in place of a union bound.

Theorem 4.4 (Uniform Convergence [Tal94]). *Let \mathcal{F} be a set of functions $f : \mathcal{X} \rightarrow \{0,1\}$ with VC dimension d . Let P be a distribution on \mathcal{X} . Then*

$$\Pr_{D \leftarrow P^n} \left[\sup_{f \in \mathcal{F}} |f(D) - f(P)| \leq \alpha \right] \geq 1 - \beta$$

whenever $n = \Omega\left(\frac{d + \log(1/\beta)}{\alpha^2}\right)$. Here $f(D) := \frac{1}{n} \sum_{x \in D} f(x)$ and $f(P) := \mathbf{E}_{X \leftarrow P}[f(X)]$.

⁴Here, for simplicity, we assume that each distribution H is given by a density function $H(\cdot)$. More generally, we define the VC dimension of \mathcal{H} to be the smallest d such that there exists a function family $\mathcal{F} \subseteq \{0,1\}^{\mathcal{X}}$ of VC dimension d with the property that, for all $H, H' \in \mathcal{H}$ we have $d_{\text{TV}}(H, H') = \sup_{f \in \mathcal{F}} \mathbf{E}_{X \leftarrow H}[f(X)] - \mathbf{E}_{X \leftarrow H'}[f(X)]$, where the supremum is over f measurable with respect to both H and H' . We ignore this technicality throughout.

It is immediate from Definition 4.2 that $VC(\mathcal{F}) \leq \lfloor \log_2 |\mathcal{F}| \rfloor$. Thus Theorem 4.4 subsumes the union bound used in the proof of Theorem 1.1.

The relevant application of uniform convergence for our algorithm is the following lemma (roughly the equivalent of Lemma 3.1), which says that good hypotheses have high scores, and bad hypotheses have low scores.

Lemma 4.5. *Let \mathcal{H} be a collection of probability distributions on \mathcal{X} with VC dimension d .*

Let $S : \mathcal{H} \times \mathcal{X}^n \rightarrow \mathbb{R}$ be as in Equation 1, namely

$$S(H, D) = \inf_{H' \in \mathcal{H}} \max \left\{ |\{x \in D : H(x) > H'(x)\}| - n \cdot (\Pr_{X \leftarrow H'}[H(X) > H'(X)] + 3\alpha), \frac{n \cdot \mathbb{I}[d_{TV}(H, H') \leq 6\alpha]}{n \cdot \mathbb{I}[d_{TV}(H, H') \leq 6\alpha]} \right\},$$

where \mathbb{I} denotes the indicator function.

Let P be a distribution on \mathcal{X} . Let $\alpha, \beta > 0$ and $n \geq O(\frac{1}{\alpha^2}(d + \log(1/\beta)))$. Suppose there exists $H^ \in \mathcal{H}$ with $d_{TV}(P, H^*) \leq \alpha$. Then, with probability at least $1 - \beta$ over $D \leftarrow P^n$, we have*

- $S(H^*, D) > \alpha n$ and
- $S(H, D) = 0$ for all $H \in \mathcal{H}$ with $d_{TV}(H, P) > 7\alpha$.

Proof. For $H, H' \in \mathcal{H}$, define $f_{H, H'} : \mathcal{X} \rightarrow \{0, 1\}$ by $f_{H, H'}(x) = 1 \iff H(x) > H'(x)$. Note that $|\{x \in D : H(x) > H'(x)\}| = \sum_{x \in D} f_{H, H'}(x)$ and d is the VC dimension of the function class $\{f_{H, H'} : H, H' \in \mathcal{H}\}$. By Theorem 4.4, if $n = \Omega\left(\frac{d + \log(1/\beta)}{\alpha^2}\right)$, then

$$\Pr_{D \leftarrow P^n} [\forall H, H' \in \mathcal{H} \quad ||\{x \in D : H(x) > H'(x)\}| - n \cdot \Pr_{X \leftarrow P}[H(X) > H'(X)]| \leq \alpha n] \geq 1 - \beta.$$

We condition on this event happening.

In order to prove the first conclusion – namely, $S(H^*, D) > \alpha n$ – it remains to show that, for all $H' \in \mathcal{H}$, we have either $d_{TV}(H^*, H') \leq 6\alpha$ or

$$|\{x \in D : H(x) > H'(x)\}| - n \cdot (\Pr_{X \leftarrow H'}[H^*(X) > H'(X)] + 3\alpha) > \alpha n.$$

If $d_{TV}(H^*, H') \leq 6\alpha$, we are done, so assume $d_{TV}(H^*, H') > 6\alpha$. By the uniform convergence event we have conditioned on,

$$\begin{aligned} |\{x \in D : H(x) > H'(x)\}| &\geq n \cdot (\Pr_{X \leftarrow P}[H(X) > H'(X)] - \alpha) \\ &\geq n \cdot (\Pr_{X \leftarrow H^*}[H(X) > H'(X)] - d_{TV}(P, H^*) - \alpha) \\ &\geq n \cdot (d_{TV}(H^*, H') + \Pr_{X \leftarrow H'}[H(X) > H'(X)] - 2\alpha) \\ &> n \cdot (6\alpha + \Pr_{X \leftarrow H'}[H(X) > H'(X)] - 2\alpha), \end{aligned}$$

from which the desired conclusion follows.

In order to prove the second conclusion – namely, $S(H, D) = 0$ for all $H \in \mathcal{H}$ with $d_{TV}(H, P) > 7\alpha$ – it suffices to show that one $H' \in \mathcal{H}$ yields a score of zero for any $H \in \mathcal{H}$ with $d_{TV}(H, P) > 7\alpha$. In particular, we show that $H' = H^*$ yields a score of zero for any such H . That is, if $d_{TV}(H, P) > 7\alpha$, then $d_{TV}(H, H^*) > 6\alpha$ and

$$|\{x \in D : H(x) > H^*(x)\}| - n \cdot (\Pr_{X \leftarrow H^*}[H(X) > H^*(X)] + 3\alpha) \leq 0.$$

By the triangle inequality $d_{TV}(H, H^*) \geq d_{TV}(H, P) - d_{TV}(P, H^*) > 7\alpha - \alpha = 6\alpha$, as required. By the uniform convergence event we have conditioned on,

$$\begin{aligned} |\{x \in D : H(x) > H^*(x)\}| &\leq n \cdot (\Pr_{X \leftarrow P}[H(X) > H^*(X)] + \alpha) \\ &\leq n \cdot (\Pr_{X \leftarrow H^*}[H(X) > H^*(X)] + d_{TV}(P, H^*) + \alpha) \\ &\leq n \cdot (\Pr_{X \leftarrow H^*}[H(X) > H^*(X)] + 2\alpha), \end{aligned}$$

which completes the proof. \square

4.2 GAP-MAX Algorithm

In place of the exponential mechanism for privately selecting a hypothesis we use the following algorithm that works under a “gap” assumption. That is, we assume that there is a $5\alpha n$ gap between the highest score and the $(k+1)$ -th highest score. Rather than paying in sample complexity for the total number of hypotheses we pay for the number of high-scoring hypotheses k .

This algorithm is based on the GAP-MAX algorithm of Bun, Dwork, Rothblum, and Steinke [BDRS18]. However, we combine their GAP-MAX algorithm with the exponential mechanism to improve the dependence on the parameter k .

Theorem 4.6. *Let \mathcal{H} and \mathcal{X} be arbitrary sets. Let $S : \mathcal{H} \times \mathcal{X}^n \rightarrow \mathbb{R}$ have sensitivity at most 1 in its second argument – that is, for all $H \in \mathcal{H}$ and all $D, D' \in \mathcal{X}^n$ differing in a single example, $|S(H, D) - S(H, D')| \leq 1$.*

For $D \in \mathcal{X}^n$ and $\alpha > 0$, define

$$K(D, 5\alpha) := \left| \left\{ H \in \mathcal{H} : S(H, D) \geq \sup_{H' \in \mathcal{H}} S(H', D) - 5\alpha n \right\} \right|.$$

Given parameters $\varepsilon, \delta, \beta > 0$ and $n, k \geq 1$, there exists a (ε, δ) -differentially private randomized algorithm $M : \mathcal{X}^n \rightarrow \mathcal{H}$ such that, for all $D \in \mathcal{X}^n$ and all $\alpha > 0$,

$$K(D, 5\alpha) \leq k \implies \Pr \left[S(M(D), D) \geq \sup_{H' \in \mathcal{H}} S(H', D) - \alpha n \right] \geq 1 - \beta$$

provided $n = \Omega \left(\frac{\min\{\log |\mathcal{H}|, \log(1/\delta)\} + \log(k/\beta)}{\alpha \varepsilon} \right)$.

Furthermore, given $\varepsilon, \beta > 0$ and $n, k \geq 1$, there exists a $\frac{1}{2}\varepsilon^2$ -concentrated differentially private [BS16] algorithm $M : \mathcal{X}^n \rightarrow \mathcal{H}$ such that, for all $D \in \mathcal{X}^n$ and all $\alpha > 0$,

$$K(D, 5\alpha) \leq k \implies \Pr \left[S(M(D), D) \geq \sup_{H' \in \mathcal{H}} S(H', D) - \alpha n \right] \geq 1 - \beta$$

provided $n = \Omega \left(\frac{\sqrt{\log |\mathcal{H}|} + \log(k/\beta)}{\alpha \varepsilon} \right)$.

Proof. We begin by describing the algorithm.

1. Let $m = \left\lceil \frac{k^2}{\beta} \right\rceil$ and let $G : \mathcal{H} \rightarrow [m]$ be a uniformly random function.⁵
2. Randomly select $B \in [m]$ with

$$\Pr[B = b] \propto \exp \left(\frac{\varepsilon}{4} \sup \{ S(H, D) : H \in \mathcal{H}, G(H) = b \} \right).$$

3. Define $\mathcal{H}_B = \{H \in \mathcal{H} : G(H) = B\}$. Let $H_B^1 = \operatorname{argmax}_{H \in \mathcal{H}_B} S(H, D)$ and $H_B^2 = \operatorname{argmax}_{H \in \mathcal{H}_B \setminus \{H_B^1\}} S(H, D)$, breaking ties arbitrarily. (That is, \mathcal{H}_B is the B -th “bin” and H_B^1 and H_B^2 are the items in this bin with the largest and second-largest scores respectively.) Define $S'_B : \mathcal{H}_B \times \mathcal{X}^n \rightarrow \mathbb{R}$ by

$$S'_B(H, D) = \frac{1}{2} \max\{0, S(H, D) - S(H_B^2, D)\}.$$

(Note that S'_B has sensitivity 1 and $S'_B(H, D) = 0$ whenever $H \neq H_B^1$.)

⁵It suffices for G to be a drawn from a universal hash function family.

4. Let \mathcal{D} be a distribution on \mathbb{R} such that adding a sample from \mathcal{D} to a sensitivity-1 function provides $(\varepsilon/4, \delta/2)$ -differential privacy (or, respectively, $\frac{1}{6}\varepsilon^2$ -concentrated differential privacy). For example, \mathcal{D} could be a Laplace distribution with scale $4/\varepsilon$ truncated to the interval $[-t, t]$ for $t = 4(1 + \log(1/\delta))/\varepsilon$ (or unbounded if $\delta = 0$). To attain concentrated differential privacy, we can set $\mathcal{D} = N(0, \frac{3}{\varepsilon^2})$, a centered Gaussian with variance $3/\varepsilon^2$.
5. Draw a sample Z_H i.i.d. from \mathcal{D} corresponding to every $H \in \mathcal{H}_B$.
6. Return $H^* = \operatorname{argmax}_{H \in \mathcal{H}_B} S'_B(H, D) + Z_H$.

The selection of B is an instantiation of the exponential mechanism [MT07] and is $(\varepsilon/2, 0)$ -differentially private. The selection of H^* in the final step is a GAP-MAX algorithm [BDRS18] and is $(\varepsilon/2, \delta)$ -differentially private. By composition, the entire algorithm is (ε, δ) -differentially private (or, respectively, $\frac{1}{2}\varepsilon^2$ -concentrated differentially private).

For the utility analysis, in order for the algorithm to output a good H^* , it suffices for the following three events to occur.

- $S(H_B^1, D) \geq \sup_{H' \in \mathcal{H}} S(H', D) - \alpha n$.
That is, restricting the search to \mathcal{H}_B , rather than all of \mathcal{H} , only reduces the score of the optimal choice by αn . The exponential mechanism ensures that this happens with probability at least $1 - \beta/4$, as long as $n \geq \frac{4 \log(2k/\beta)}{\varepsilon \alpha}$.
- $S(H_B^2, D) < \sup_{H' \in \mathcal{H}} S(H', D) - 5\alpha n$.
That is, the second-highest score within \mathcal{H}_B is at least $5\alpha n$ less than the highest score overall. We have assumed that there are at most k elements $H \in \mathcal{H}$ such that $S(H, D) \geq \sup_{H' \in \mathcal{H}} S(H', D) - 5\alpha n$. Call these “large elements.” Since G is random and $m \geq k^2/\beta$, the probability that more than one large element satisfies $G(H) = B$ is at most $\beta/2$. That is to say, with high probability there are no collisions under the hash function G of the k large elements. This suffices for the event to occur.
- $\sup_{H \in \mathcal{H}_B} |Z_H| \leq \alpha n$.
If the noise distribution \mathcal{D} is supported on $[-\alpha n, \alpha n]$, then this condition holds with probability 1. For the truncated Laplace distribution, this is possible whenever $n \geq 1 + 4 \log(1/\delta)/\alpha \varepsilon$. Alternatively, we can use unbounded Laplace noise and a union bound to show that this event occurs with probability at least $1 - \beta/4$ whenever $n \geq 4 \log(4|\mathcal{H}_B|/\beta)/\varepsilon \alpha$. For Gaussian noise, $n \geq \frac{3}{\varepsilon \alpha} \sqrt{\log(4|\mathcal{H}_B|/\beta)}$ suffices.

Assuming the first and second events occur, we have $S'_B(H_B^1, D) = \frac{S(H_B^1, D) - S(H_B^2, D)}{2} > 2\alpha n$. Given this, the third event implies $H^* = H_B^1$. Finally, the first event then implies $S(H^*, D) \geq \sup_{H' \in \mathcal{H}} S(H', D) - \alpha n$, as required. A union bound over the three events completes the proof. \square

Now we can combine the VC-based uniform convergence bound with the GAP-MAX algorithm to prove our result.

Proof of Theorem 4.1. By Lemma 4.5, with high probability over the draw of the dataset D , our score function satisfies $\sup_{H \in \mathcal{H}} S(H, D) \geq S(H^*, D) > \alpha n$ and $S(H, D) = 0$ whenever $d_{\text{TV}}(H, P) > 7\alpha$. This requires $n = \Omega(d/\alpha^2)$.

Note that the score function S has sensitivity-1, since it is the supremum of counts. Conditioned on the uniform convergence event, the maximum score is at least αn and there are at most k elements of \mathcal{H} with score greater than 0. Thus we can apply the GAP-MAX algorithm of Theorem 4.6. If $n = \Omega((\min\{\log |\mathcal{H}|, \log(1/\delta)\} + \log(k))/\alpha \varepsilon)$, then with high probability, the algorithm outputs $\hat{H} \in \mathcal{H}$ with score at least $\frac{4}{5}\alpha n$, as required. \square

5 Packings, Lower Bounds, and Relations to Covers

In this section, we show that the sample complexity of our algorithms for private hypothesis selection with pure differential privacy cannot be improved, at least for constant values of the proximity parameter α . We first apply a packing argument [HT10, BBKN14] to show a lower bound which is logarithmic in the packing number of the class of distributions (Lemma 5.1). We then state a folklore relationship between the sizes of maximal packings and minimal covers (Lemma 5.2), which shows that instantiating our private hypothesis selection algorithm with a minimal cover gives essentially optimal sample complexity (Theorem 5.3).

Lemma 5.1. *Suppose there exists an α -packing \mathcal{P}_α of a set of distributions \mathcal{H} . Then any ε -differentially private algorithm which takes as input samples $X_1, \dots, X_n \sim P$ for some $P \in \mathcal{H}$ and produces a distribution \hat{H} such that $d_{\text{TV}}(P, \hat{H}) \leq \alpha$ with probability $\geq 9/10$ requires*

$$n = \Omega\left(\frac{\log |\mathcal{P}_\alpha|}{\varepsilon}\right).$$

Proof. Let M be a ε -differentially private algorithm with the stated accuracy requirement, and denote by $M(P^n)$ the distribution on hypotheses obtained by running M on n i.i.d. samples from a distribution $P \in \mathcal{H}$. For each $P \in \mathcal{P}_\alpha$, let B_P denote the set of distributions which are at total variation distance at most α from P . Then the accuracy requirement implies that $\Pr_{\hat{H} \leftarrow M(P^n)}[\hat{H} \in B_P] \geq 9/10$. Let $P_0 \in \mathcal{P}_\alpha$ be an arbitrary packing element. Then by group privacy applied to groups of size n , we have

$$\Pr_{\hat{H} \leftarrow M(P_0^n)}[\hat{H} \in B_P] \geq e^{-\varepsilon n} \cdot 9/10$$

for every $P \in \mathcal{P}_\alpha$. The fact that \mathcal{P}_α is an α -packing implies that the sets B_P are all disjoint, and hence

$$1 \geq \sum_{P \in \mathcal{P}_\alpha} \Pr_{\hat{H} \leftarrow M(P_0^n)}[\hat{H} \in B_P] \geq |\mathcal{P}_\alpha| \cdot e^{-\varepsilon n} \cdot 9/10.$$

Rearranging gives us the stated lower bound on n . \square

Lemma 5.2. *For a set of distributions \mathcal{H} , let p_α and c_α be the size of the largest α -packing and smallest α -cover of \mathcal{H} , respectively. Then*

$$p_{2\alpha} \leq c_\alpha \leq p_\alpha.$$

Proof. We first prove the inequality on the left. Let \mathcal{C}_α be a cover of \mathcal{H} of size c_α . If $c_\alpha = \infty$, we are done. Otherwise, let S be any set of points of size at least $c_\alpha + 1$. By the pigeonhole principle, there exists $P \in \mathcal{C}_\alpha$ and two distributions $Q, Q' \in S$ such that $d_{\text{TV}}(P, Q) \leq \alpha$ and $d_{\text{TV}}(P, Q') \leq \alpha$. Hence $d_{\text{TV}}(Q, Q') \leq 2\alpha$ by the triangle inequality, so S cannot be (2α) -packing of \mathcal{H} . This suffices to show that $p_{2\alpha} \leq c_\alpha$.

Next, we prove the inequality on the right. Let \mathcal{P}_α be a maximal α -packing with size $|\mathcal{P}_\alpha| = p_\alpha$. If $p_\alpha = \infty$, we are done. Otherwise, we claim that \mathcal{P}_α is also an α -cover of \mathcal{H} , and hence $c_\alpha \leq |\mathcal{P}_\alpha| = p_\alpha$. To see this, suppose for the sake of contradiction that there were a distribution $P \in \mathcal{H}$ with $d_{\text{TV}}(P, \mathcal{P}_\alpha) > \alpha$. Then we could add P to \mathcal{P}_α to produce a strictly larger packing, contradicting the maximality of \mathcal{P}_α . \square

Theorem 5.3. *Let \mathcal{H} be a set of distributions, and let n_α^* denote the minimum number of samples such that there exists an ε -differentially private algorithm which takes as input samples $X_1, \dots, X_{n_\alpha^*} \sim P$ for some $P \in \mathcal{H}$ and outputs a distribution \hat{H} such that $d_{\text{TV}}(P, \hat{H}) \leq \alpha$ with probability $\geq 9/10$. Then there exists a cover of \mathcal{H} such that the instantiation of the algorithm underlying Theorem 1.1 with this cover outputs a \hat{H} such that $d_{\text{TV}}(P, \hat{H}) \leq 7\alpha$ with probability $\geq 9/10$ for any $n = \Omega(n_\alpha^* \cdot (\varepsilon/\alpha^2 + 1/\alpha))$.*

Proof. Let p_α denote the size of the largest α -packing of \mathcal{H} . By Lemma 5.1, we have $n_\alpha^* = \Omega(\log p_\alpha / \varepsilon)$. On the other hand, by Lemma 5.2, we know that there exists an α -cover \mathcal{C}_α of \mathcal{H} with $|\mathcal{C}_\alpha| \leq p_\alpha$. Hence $\log |\mathcal{C}_\alpha| \leq O(\varepsilon \cdot n_\alpha^*)$ and the asserted sample complexity guarantee follows from Corollary 1.2. \square

6 Applications of Hypothesis Selection

In this section, we give a number of applications of Theorem 1.1, primarily to obtain sample complexity bounds for learning a number of distribution classes of interest. Recall Corollary 1.2, which is an immediate corollary of Theorem 1.1. This indicates that we can privately semi-agnostically learn a class of distributions with a number of samples proportional to the logarithm of its covering number.

Corollary 1.2. *Suppose there exists an α -cover \mathcal{C}_α of a set of distributions \mathcal{H} , and that we are given a set of samples $X_1, \dots, X_n \sim P$, where $d_{\text{TV}}(P, \mathcal{H}) \leq \alpha$. For any constant $\zeta > 0$, there exists an ε -differentially private algorithm (with respect to the input $\{X_1, \dots, X_n\}$) which outputs a distribution $H^* \in \mathcal{C}_\alpha$ such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, as long as*

$$n = \Omega\left(\frac{\log |\mathcal{C}_\alpha|}{\alpha^2} + \frac{\log |\mathcal{C}_\alpha|}{\alpha\varepsilon}\right).$$

Note that the factor of $(6 + 2\zeta)\alpha$ in the corollary statement (versus $(3 + \zeta)\alpha$ in the statement of Theorem 1.1) is due to the fact the algorithm is semi-agnostic, and the closest element in the cover is 2α -close to P , rather than just α -close.

We instantiate this result to give the sample complexity results for semi-agnostically learning product distributions (Section 6.1), Gaussian distributions (Section 6.2), sums of some independent random variable classes (Section 6.3), piecewise polynomials (Section 6.4), and mixtures (Section 6.5). Furthermore, we mention an application to private PAC learning (Section 6.6), when the distribution of unlabeled examples is known to come from some hypothesis class.

6.1 Product Distributions

As a first application, we first give an ε -differentially private algorithm for learning product distributions over discrete alphabets.

Definition 6.1. *A (k, d) -product distribution is a distribution over $[k]^d$, such that its marginal distributions are independent (i.e., the distribution is the product of its marginals).*

We start by constructing a cover for product distributions.

Lemma 6.2. *There exists an α -cover of the set of (k, d) -product distributions of size*

$$O\left(\frac{kd}{\alpha}\right)^{d(k-1)}.$$

Proof. Consider some fixed product distribution P , with marginal distributions (P_1, \dots, P_d) . We will construct a cover that contains a distribution Q (with marginals (Q_1, \dots, Q_d)) that is α -close in total variation distance.

First, by triangle inequality, we have that $d_{\text{TV}}(P, Q) \leq \sum_{i=1}^d d_{\text{TV}}(P_i, Q_i)$, so it suffices to approximate each marginal distribution to accuracy α/d . Stated another way, we must generate an (α/d) -cover of distributions over $[k]$, and we can then take its d -wise Cartesian product. Raising the size of this underlying cover to the power d gives us the size of the overall cover.

To (α/d) -cover a distribution over $[k]$, we will additively grid the probability of each symbol at granularity $\Theta(\frac{\alpha}{kd})$, choosing the probability of the last symbol k such that the sum is normalized. This will incur $\Theta(\frac{\alpha}{kd})$ error per symbol (besides for symbol k), and summing over the $k-1$ symbols accumulates error $\Theta(\frac{\alpha}{d})$. It can also be argued that the error on symbol k is $O(\frac{\alpha}{d})$ – with an appropriate choice of granularity, this gives us an (α/d) -cover for distributions over $[k]$. The size of this cover is $O(\frac{kd}{\alpha})^{k-1}$, which allows us to conclude the lemma statement. \square

With this cover in hand, applying Corollary 1.2 allows us to conclude the following sample complexity upper bound.

Corollary 6.3. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is α -close to a (k, d) -product distribution. Then for any constant $\zeta > 0$, there exists an ε -differentially private algorithm which outputs a (k, d) -product distribution H^* such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, so long as*

$$n = \Omega\left(kd \log\left(\frac{kd}{\alpha}\right) \left(\frac{1}{\alpha^2} + \frac{1}{\alpha\varepsilon}\right)\right).$$

This gives the first $\tilde{O}(d)$ sample algorithm for learning a binary product distribution in total variation distance under pure differential privacy, improving upon the work of Kamath, Li, Singhal, and Ullman [KLSU19] by strengthening the privacy guarantee at a minimal cost in the sample complexity. The natural way to adapt their result from concentrated to pure differential privacy would require $\Omega(d^{3/2})$ samples.

Remark 6.4. *Properly learning a product distribution over $\{0, 1\}^d$ to total variation distance $\leq \frac{1}{2}$ implies learning its mean $\mu \in [0, 1]^d$ up to ℓ_1 error $\leq 2\sqrt{d}$; see Lemma 6.5 below.*

Thus Corollary 6.3 implies a ε -differentially private algorithm which takes $n = \tilde{O}(d/\varepsilon)$ samples from a product distribution P on $\{0, 1\}^d$ and, with high probability, outputs an estimate $\hat{\mu}$ of its mean μ with $\|\hat{\mu} - \mu\|_1 \leq 2\sqrt{d}$.

In contrast, for non-product distributions over the hypercube, estimating the mean to the same accuracy under ε -differential privacy requires $n = \Omega(d^{3/2}/\varepsilon)$ samples [HT10, SU15]. Thus we have a polynomial separation between estimating product and non-product distributions under pure differential privacy.

Lemma 6.5. *If P and Q are product distributions on \mathbb{R}^d with $d_{\text{TV}}(P, Q) \leq \frac{1}{2}$ and per-coordinate variance at most σ^2 , then*

$$\|\mathbf{E}_{X \leftarrow P}[X] - \mathbf{E}_{X \leftarrow Q}[X]\|_1 \leq 4\sqrt{d\sigma^2}.$$

Proof. Let $\mu = \mathbf{E}_{X \leftarrow P}[X] \in \mathbb{R}^d$ and $\mu' = \mathbf{E}_{X \leftarrow Q}[X] \in \mathbb{R}^d$. Let $\tau = \|\mu - \mu'\|_1$. Let $\nu = \text{sign}(\mu - \mu') \in$

$\{-1, +1\}^d$ so that $\langle \nu, \mu - \mu' \rangle = \tau$. We have

$$\begin{aligned}
\frac{1}{2} &\geq d_{\text{TV}}(P, Q) \geq \Pr_{X \leftarrow P}[\langle \nu, X \rangle \geq t] - \Pr_{X \leftarrow Q}[\langle \nu, X \rangle \geq t] \\
&= \Pr_{X \leftarrow P}[\langle \nu, X - \mu \rangle \geq t - \langle \nu, \mu \rangle] - \Pr_{X \leftarrow Q}[\langle \nu, X - \mu' \rangle \geq t - \langle \nu, \mu \rangle + \langle \nu, \mu - \mu' \rangle] \\
(\text{set } t = \langle \nu, \mu \rangle - \frac{\tau}{2}) \quad &= \Pr_{X \leftarrow P}[\langle \nu, X - \mu \rangle \geq -\frac{\tau}{2}] - \Pr_{X \leftarrow Q}[\langle \nu, X - \mu' \rangle \geq +\frac{\tau}{2}] \\
&= 1 - \Pr_{X \leftarrow P}[\langle \nu, X - \mu \rangle < -\frac{\tau}{2}] - \Pr_{X \leftarrow Q}[\langle \nu, X - \mu' \rangle \geq +\frac{\tau}{2}] \\
(\text{Chebyshev's inequality}) \quad &\geq 1 - \frac{\mathbf{E}_{X \leftarrow P}[\langle \nu, X - \mu \rangle^2]}{(\tau/2)^2} - \frac{\mathbf{E}_{X \leftarrow Q}[\langle \nu, X - \mu' \rangle^2]}{(\tau/2)^2} \\
&= 1 - \frac{4}{\tau^2} \sum_{i=1}^d \mathbf{E}_{X \leftarrow P}[(X_i - \mu_i)^2] + \mathbf{E}_{X \leftarrow Q}[(X_i - \mu'_i)^2] \\
&\geq 1 - \frac{8d\sigma^2}{\tau^2}.
\end{aligned}$$

Rearranging yields $\tau \leq 4\sqrt{d\sigma^2}$, as required. \square

6.2 Gaussian Distributions

We next give private algorithms for learning Gaussian distributions.

Definition 6.6. A Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ in \mathbb{R}^d is a distribution with PDF

$$p(x) = \frac{\exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)}{\sqrt{(2\pi)^d |\Sigma|}}.$$

We describe covers for Gaussian distributions with known and unknown covariance.

Lemma 6.7. *There exists an α -cover of the set of Gaussian distributions $\mathcal{N}(\mu, I)$ in d dimensions with $\|\mu\|_2 \leq R$ of size*

$$O\left(\frac{dR}{\alpha}\right)^d.$$

Proof. It is well-known that estimating a Gaussian distribution with unknown mean in total variation distance corresponds to estimating μ in ℓ_2 -distance (see, e.g., [DKK⁺16]). By the triangle inequality, in order to α -cover the space, it suffices to (α/d) -cover each standard basis direction. Since we know the mean in each direction is bounded by R , a simple additive grid in each direction with granularity $\Theta\left(\frac{\alpha}{d}\right)$ will suffice, resulting in a cover for each direction of size $O\left(\frac{dR}{\alpha}\right)$. Taking the Cartesian product over d dimensions gives the desired result. \square

Lemma 6.8. *There exists an α -cover of the set of Gaussian distributions $\mathcal{N}(\mu, \Sigma)$ in d -dimensions with $\|\mu\|_2 \leq R$ and $I \preceq \Sigma \preceq \kappa I$ of size*

$$O\left(\frac{dR}{\alpha}\right)^d \cdot O\left(\frac{d\kappa}{\alpha}\right)^{d(d+1)/2}.$$

Proof. The former term is obtained similarly to the expression in Lemma 6.7. Since $I \preceq \Sigma$, we can still bound the total variation contribution by the ℓ_2 -distance between the mean vectors. We thus turn our attention to the latter term. To construct our cover, we must argue about the total

variation distance between $\mathcal{N}(0, \Sigma)$ and $\mathcal{N}(0, \hat{\Sigma})$. If $|\Sigma(i, j) - \hat{\Sigma}(i, j)| \leq \gamma$, and $I \preceq \Sigma$, Proposition 32 of [VV10] implies:

$$d_{\text{TV}}(\mathcal{N}(0, \Sigma), \mathcal{N}(0, \hat{\Sigma})) \leq O(d\gamma).$$

We will thus perform a gridding, in order to approximate each entry of Σ to an additive $O(\gamma) = O(\alpha/d)$. However, in order to ensure that the resulting matrix is PSD, we grid over entries of $\hat{\Sigma}$'s Cholesky decomposition, rather than grid for $\hat{\Sigma}$ itself. Since the largest element of Σ is bounded by κ , the largest element of its Cholesky decomposition must be bounded by $\sqrt{\kappa}$. An additive grid over the range $[0, \sqrt{\kappa}]$ with granularity $O(\gamma/\sqrt{\kappa})$ suffices to get $\hat{\Sigma}$ which bounds the entrywise distance as $O(\gamma)$. This requires $O(d\kappa/\alpha)$ candidates per entry, and we take the Cartesian product over all $d(d+1)/2$ entries of the Cholesky decomposition, giving the desired result. \square

In addition, we can obtain bounds of the VC dimension of the Scheffé sets of Gaussian distributions.

Lemma 6.9. *The set of Gaussian distributions with fixed variance – i.e., all $\mathcal{N}(\mu, I)$ with $\mu \in \mathbb{R}^d$ – has VC dimension $d+1$. Furthermore, the set of Gaussians with unknown variance – i.e., all $\mathcal{N}(\mu, \Sigma)$ with $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathbb{R}^{d \times d}$ positive definite – has VC dimension $O(d^2)$.*

Proof. For Gaussians with fixed variance, the Scheffé sets correspond to linear threshold functions, which have VC dimension $d+1$. For Gaussians with unknown variance, the Scheffé sets correspond to quadratic threshold functions, which have VC dimension $\binom{d+2}{2} = O(d^2)$ [Ant95]. \square

Combining the covers of Lemmas 6.7 and 6.8 and the VC bound of Lemma 6.9 with Theorem 4.1 implies the following corollaries for Gaussian estimation.

Corollary 6.10. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is α -close to a Gaussian distribution $\mathcal{N}(\mu, I)$ in d -dimensions with $\|\mu\| \leq R$. Then for any constant $\zeta > 0$, there exists an ε -differentially private algorithm which outputs a Gaussian distribution H^* such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, so long as*

$$n = \Omega\left(\frac{d}{\alpha^2} + \frac{d}{\alpha\varepsilon} \log\left(\frac{dR}{\alpha}\right)\right).$$

Corollary 6.11. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is α -close to a Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ in d -dimensions with $\|\mu\| \leq R$ and $I \preceq \Sigma \preceq \kappa I$. Then for any constant $\zeta > 0$, there exists an ε -differentially private algorithm which outputs a Gaussian distribution H^* such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, so long as*

$$n = \Omega\left(\frac{d^2}{\alpha^2} + \frac{1}{\alpha\varepsilon} \left(d \log\left(\frac{dR}{\alpha}\right) + d^2 \log\left(\frac{d\kappa}{\alpha}\right)\right)\right).$$

Similar to the product distribution case, these are the first $\tilde{O}(d)$ and $\tilde{O}(d^2)$ sample algorithms for learning Gaussians total variation distance under pure differential privacy, improving upon the concentrated differential privacy results of Kamath, Li, Singhal, and Ullman [KLSU19].

6.2.1 Gaussians with Unbounded Mean

Extending Corollary 6.10, we consider multivariate Gaussian hypotheses with known covariance and unknown mean, *without* assuming bound on the mean (the parameter R in the discussion above). To handle the unbounded mean we must relax to approximate differential privacy.

In place of Lemma 6.7, we construct a locally small cover:

Lemma 6.12. *For any $d \in \mathbb{N}$ and $\alpha \in (0, 1/30]$, there exists an α -cover \mathcal{C}_α of the set of Gaussian distributions $\mathcal{N}(\mu, I)$ in d dimensions satisfying*

$$\forall \mu \in \mathbb{R}^d \quad |\{H \in \mathcal{C}_\alpha : d_{\text{TV}}(H, \mathcal{N}(\mu, I)) \leq 7\alpha\}| \leq 2^{15d}.$$

Proof. For $\mu, \mu' \in \mathbb{R}^d$, we have

$$\begin{aligned} d_{\text{TV}}(\mathcal{N}(\mu, I), \mathcal{N}(\mu', I)) &= 2\Pr\left[\mathcal{N}(0, 1) \in \left[0, \frac{1}{2}\|\mu - \mu'\|_2\right]\right] \\ &= \sqrt{\frac{2}{\pi}} \int_0^{\frac{1}{2}\|\mu - \mu'\|_2} e^{-x^2/2} dx \\ &\leq \frac{\|\mu - \mu'\|_2}{\sqrt{2\pi}}. \end{aligned}$$

Furthermore, for any $c > 0$,

$$d_{\text{TV}}(\mathcal{N}(\mu, I), \mathcal{N}(\mu', I)) \geq \begin{cases} \frac{\|\mu - \mu'\|_2}{\sqrt{2\pi}} \cdot e^{-c^2/2} & \text{if } \frac{1}{2}\|\mu - \mu'\|_2 \leq c \\ \frac{c \cdot e^{-c^2/2}}{\sqrt{2\pi}} & \text{if } \frac{1}{2}\|\mu - \mu'\|_2 \geq c \end{cases}.$$

Let

$$\mathcal{C}_\alpha = \left\{ \mathcal{N}\left(m \cdot \frac{\alpha\sqrt{8\pi}}{\sqrt{d}}, I\right) : m \in \mathbb{Z}^d \right\}.$$

Fix $\mu \in \mathbb{R}^d$. Let $\mu^* = \mu \frac{\sqrt{d}}{\alpha\sqrt{8\pi}} \in \mathbb{R}^d$ and let $m \in \mathbb{Z}^d$ be μ^* rounded to the nearest integer coordinate-wise, so that $\|m - \mu^*\|_\infty \leq \frac{1}{2}$. Then

$$\begin{aligned} d_{\text{TV}}\left(\mathcal{N}(\mu, I), \mathcal{N}\left(m \cdot \frac{\alpha\sqrt{8\pi}}{\sqrt{d}}, I\right)\right) &= d_{\text{TV}}\left(\mathcal{N}\left(\mu^* \cdot \frac{\alpha\sqrt{8\pi}}{\sqrt{d}}, I\right), \mathcal{N}\left(m \cdot \frac{\alpha\sqrt{8\pi}}{\sqrt{d}}, I\right)\right) \\ &\leq \frac{1}{\sqrt{2\pi}} \frac{\alpha\sqrt{8\pi}}{\sqrt{d}} \|\mu^* - m\|_2 \\ &\leq \alpha, \end{aligned}$$

since $\|\mu^* - m\|_2 \leq \sqrt{d}\|\mu^* - m\|_\infty \leq \frac{\sqrt{d}}{2}$. This proves that \mathcal{C}_α is a α -cover of $\{\mathcal{N}(\mu, I) : \mu \in \mathbb{R}^d\}$.

It remains to show that the cover is “locally small”. Let $m' \in \mathbb{Z}^d$. Then

$$\begin{aligned} d_{\text{TV}}\left(\mathcal{N}(\mu, I), \mathcal{N}\left(m' \cdot \frac{\alpha\sqrt{8\pi}}{\sqrt{d}}, I\right)\right) &= d_{\text{TV}}\left(\mathcal{N}\left(\mu^* \cdot \frac{\alpha\sqrt{8\pi}}{\sqrt{d}}, I\right), \mathcal{N}\left(m' \cdot \frac{\alpha\sqrt{8\pi}}{\sqrt{d}}, I\right)\right) \\ &\geq \frac{c \cdot e^{-c^2/2}}{\sqrt{2\pi}} \quad \text{if } \frac{1}{2}\|\mu^* - m'\|_2 \frac{\alpha\sqrt{8\pi}}{\sqrt{d}} \geq c \\ &> 7\alpha \quad \text{if } \|\mu^* - m'\|_2 \geq 30 \frac{\sqrt{d}}{\sqrt{2\pi}}, \end{aligned}$$

where the final inequality follows by setting $c = 30\alpha \leq 1$. Thus

$$\begin{aligned}
|\{H \in \mathcal{C}_\alpha : d_{\text{TV}}(H, \mathcal{N}(\mu, I)) \leq 7\alpha\}| &\leq \left| \left\{ m' \in \mathbb{Z}^d : \|\mu^* - m'\|_2 < 30 \frac{\sqrt{d}}{\sqrt{2\pi}} \right\} \right| \\
&\leq \left| \left\{ m' \in \mathbb{Z}^d : \|m - m'\|_2 < 30 \frac{\sqrt{d}}{\sqrt{2\pi}} + \|\mu^* - m'\|_2 \right\} \right| \\
&\leq \left| \left\{ m' \in \mathbb{Z}^d : \|m - m'\|_2 < 13\sqrt{d} \right\} \right| \\
&\leq \left| \left\{ w \in \mathbb{Z}^d : \|w\|_1 < 13d \right\} \right|.
\end{aligned}$$

Now we note that any $w \in \mathbb{Z}^d$ with $\|w\|_1 \leq r$ can be written as $w = x - y$ where $x, y \in \mathbb{Z}^d$ with $\sum_{i=1}^d x_i + y_i = r$ and, for all $i \in [d]$, we have $x_i \geq 0$ and $y_i \geq 0$. Instead of counting these w vectors, we can count such (x, y) vector pairs. We can interpret a pair of x, y vectors as a way of putting r balls into $2d$ bins or r “stars” and $2d - 1$ “bars”. We can thus count

$$\left| \left\{ w \in \mathbb{Z}^d : \|w\|_1 < 13d \right\} \right| \leq \left| \left\{ x, y \in \mathbb{Z}^d : \|x\|_1 + \|y\|_2 = 13d - 1, x \geq 0, y \geq 0 \right\} \right| \leq \binom{15d - 2}{2d - 1} \leq 2^{15d}.$$

□

Applying Theorem 4.1 with the cover of Lemma 6.12 and the VC bound from Lemma 6.9 now yields an algorithm.

Corollary 6.13. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is a spherical Gaussian distribution $\mathcal{N}(\mu, I)$ in d -dimensions. Then there exists a (ε, δ) -differentially private algorithm which outputs a spherical Gaussian distribution H^* such that $d_{\text{TV}}(P, H^*) \leq 7\alpha$ with probability $\geq 1 - 2^{-d}$, so long as*

$$n = \Omega \left(\frac{d}{\alpha^2} + \frac{d + \log(1/\delta)}{\alpha\varepsilon} \right).$$

Karwa and Vadhan [KV18] give an algorithm for estimating a univariate Gaussian with unbounded mean. One can consider applying their algorithm independently to the d coordinates (which is done in [KLSU19]), giving a sample complexity bound of $\tilde{O} \left(\frac{d}{\alpha^2} + \frac{d}{\alpha\varepsilon} + \frac{\sqrt{d} \log^{3/2}(1/\delta)}{\varepsilon} \right)$, which our bound dominates except for very small values of α .

6.2.2 Univariate Gaussians with Unbounded Mean and Variance

Our methods also allow us to derive learning algorithms for univariate Gaussians with unknown mean and variance.

Lemma 6.14. *For all α less than some absolute constant, there exists an α -cover \mathcal{C}_α of the set of univariate Gaussian distributions satisfying*

$$\forall \mu, \sigma \in \mathbb{R} \quad \left| \{H \in \mathcal{C}_\alpha : d_{\text{TV}}(H, \mathcal{N}(\mu, \sigma^2)) \leq 7\alpha\} \right| \leq O(1).$$

Proof. For all $\mu, \tilde{\mu} \in \mathbb{R}$ and all $\sigma, \tilde{\sigma} > 0$, we have [DMR18, Thm 1.3]

$$\frac{1}{200} \min \left\{ 1, \max \left\{ \frac{|\tilde{\sigma}^2 - \sigma^2|}{\tilde{\sigma}^2}, \frac{40|\tilde{\mu} - \mu|}{\tilde{\sigma}} \right\} \right\} \leq d_{\text{TV}}(\mathcal{N}(\mu, \sigma^2), \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2)) \leq \frac{3|\tilde{\sigma}^2 - \sigma^2|}{2\tilde{\sigma}^2} + \frac{|\tilde{\mu} - \mu|}{2\tilde{\sigma}}.$$

Let $\beta = \alpha$ and $\gamma = \log(1 + \alpha/2)$. Define the set of distributions

$$\mathcal{C}_\alpha = \{\mathcal{N}(\beta e^{\gamma n} m, e^{2\gamma n}) : n, m \in \mathbb{Z}\}.$$

We first show that \mathcal{C}_α is an α -cover: Let $\mu \in \mathbb{R}$ and $\sigma > 0$. Let $n = \left\lceil \frac{\log \sigma}{\gamma} \right\rceil$ and $m = \left\lceil \frac{\mu}{\beta e^{\gamma n}} \right\rceil$, where $[x]$ denotes the nearest integer to x , satisfying $|x - [x]| \leq \frac{1}{2}$. Let $\tilde{\sigma} = e^{\gamma n}$ and $\tilde{\mu} = \beta e^{\gamma n} m$ so that $e^{-\gamma} \leq \frac{\tilde{\sigma}^2}{\sigma^2} \leq e^\gamma$ and $|\mu - \tilde{\mu}| \leq \frac{1}{2} \beta e^{\gamma n} = \frac{1}{2} \beta \tilde{\sigma}$. Thus $\mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2) \in \mathcal{C}_\alpha$ and $d_{\text{TV}}(\mathcal{N}(\mu, \sigma^2), \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2)) \leq \frac{3}{2}(e^\gamma - 1) + \frac{\beta}{4} \leq \alpha$, as required.

It only remains to show that the cover size is locally small. Let $\mu \in \mathbb{R}$ and $\sigma > 0$.

$$\begin{aligned} |\{H \in \mathcal{C}_\alpha : d_{\text{TV}}(H, \mathcal{N}(\mu, \sigma^2)) \leq 7\alpha\}| &= |\{n, m \in \mathbb{Z} : d_{\text{TV}}(\mathcal{N}(\beta e^{\gamma n} m, e^{2\gamma n}), \mathcal{N}(\mu, \sigma^2)) \leq 7\alpha\}| \\ &\leq \left| \left\{ n, m \in \mathbb{Z} : \max \left\{ \frac{|e^{2\gamma n} - \sigma^2|}{e^{2\gamma n}}, \frac{40|\beta e^{\gamma n} m - \mu|}{e^{\gamma n}} \right\} \leq 1400\alpha \right\} \right| \\ &= \left| \left\{ n, m \in \mathbb{Z} : \begin{array}{l} \frac{-\log(1+1400\alpha)}{2\gamma} \leq n - \frac{\log \sigma}{\gamma} \leq \frac{-\log(1-1400\alpha)}{2\gamma} \\ -35\frac{\alpha}{\beta} \leq m - \frac{\mu}{\beta e^{\gamma n}} \leq 35\frac{\alpha}{\beta} \end{array} \right\} \right| \\ &\leq \left(\frac{-\log(1-1400\alpha)}{2\gamma} - \frac{-\log(1+1400\alpha)}{2\gamma} + 1 \right) \cdot (35 - (-35) + 1) \\ &= \frac{1}{2\log(1+\alpha/2)} \log \left(\frac{1+1400\alpha}{1-1400\alpha} \right) \cdot 71 + 71 \\ &= O(1). \end{aligned}$$

□

Combining Lemma 6.14 with Lemma 6.9 and Theorem 4.1 yields the following.

Corollary 6.15. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is a univariate Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$. Then there exists a (ε, δ) -differentially private algorithm which outputs a univariate Gaussian distribution H^* such that $d_{\text{TV}}(P, H^*) \leq 7\alpha$ with probability $\geq 9/10$, so long as*

$$n = \Omega \left(\frac{1}{\alpha^2} + \frac{\log(1/\delta)}{\alpha\varepsilon} \right).$$

This sample complexity is comparable to that of Karwa and Vadhan [KV18], who give an (ε, δ) -DP algorithm with sample complexity $\tilde{O} \left(\frac{1}{\alpha^2} + \frac{1}{\alpha\varepsilon} + \frac{\log(1/\delta)}{\varepsilon} \right)$.

6.3 Sums of Independent Random Variables

In this section, we apply our results to distribution classes which are defined as the sum of independent (but not necessarily identical) distributions. These are all generalizations of the classical Binomial distribution, and they have enjoyed a great deal of study into the construction of sparse covers. To the best of our knowledge, we are the first to provide private learning algorithms for these classes.

We start with the Poisson Binomial distribution.

Definition 6.16. *A k -Poisson Binomial Distribution (k -PBD) is the sum of k independent Bernoulli random variables.*

We next consider sums of independent integer random variables, which generalize PBDs (which correspond to the case $d = 2$).

Definition 6.17. A (k, d) -Sum of Independent Integer Random Variables $((k, d)$ -SIIRV) is the sum of k independent random variables over $\{0, \dots, d-1\}$.

Finally, we consider Poisson Multinomial distributions, which again generalize PBDs (which, again, correspond to the case $d = 2$).

Definition 6.18. A (k, d) -Poisson Multinomial Distribution $((k, d)$ -PMD) is the sum of k independent d -dimensional categorical random variables, i.e., distributions over $\{e_1, \dots, e_d\}$, where e_i is the i th basis vector.

We start with a covering result for SIIRVs (including the special case of PBDs), which appears in [DKS16b]. Previous covers for PBDs and SIIRVs appear in [DP09, DP15b, DDO⁺13].

Lemma 6.19 ([DKS16b]). *There exists an α -cover of the set of (k, d) -SIIRVs of size*

$$k \cdot 2^{O(d \log^2(1/\alpha) + d \log^2 d)}.$$

Using this cover, we can apply Corollary 1.2 to attain the following learning result for PBDs and SIIRVs.

Corollary 6.20. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is α -close to a (k, d) -SIIRV. Then for any constant $\zeta > 0$, there exists an ε -differentially private algorithm which outputs a (k, d) -SIIRV H^* such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, so long as*

$$n = \Omega \left((\log k + d \log^2(1/\alpha) + d \log^2 d) \left(\frac{1}{\alpha^2} + \frac{1}{\alpha \varepsilon} \right) \right).$$

Next, we move on to PMDs. The following cover does not appear verbatim in any single location, but is a combination of results from a few different sources. The proofs for the best bounds on first term appears in [DDKT16], the second in [DKT15], and the third in [DKS16a]. Larger covers previously appeared in [DP08, DP15a].

Lemma 6.21 ([DKT15, DDKT16, DKS16a]). *For any $d > 2$, there exists an α -cover of the set of (k, d) -PMDs of size*

$$k^{O(d)} \cdot \min \left\{ 2^{\text{poly}(d/\alpha)}, (1/\alpha)^{O(d \log(d/\alpha) / \log \log(d/\alpha))^{d-1}} \right\}.$$

This implies the following learning result for PMDs.

Corollary 6.22. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is α -close to a (k, d) -PMD, for any $d > 2$. Then there exists an ε -differentially private algorithm which outputs a (k, d) -PMD H^* such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, so long as*

$$n = \tilde{\Omega} \left(\left(d \log k + \min \left\{ \text{poly} \left(\frac{d}{\alpha} \right), O \left(\frac{d \log(d/\alpha)}{\log \log(d/\alpha)} \right)^{d-1} \cdot \log(1/\alpha) \right\} \right) \left(\frac{1}{\alpha^2} + \frac{1}{\alpha \varepsilon} \right) \right).$$

6.4 Piecewise Polynomials

In this section, we apply our results to semi-agnostically learn piecewise polynomials. This class of distributions is very expressive, allowing us to approximate a wide range of natural distribution classes.

Definition 6.23. A (t, d, k) -piecewise polynomial distribution is a distribution P over $[k]$, such that there exists a partition of $[k]$ into t disjoint intervals I_1, \dots, I_t such that on each interval $I_j \subseteq [k]$, the probability mass function of P takes the form $p_j(x) = \sum_{i=0}^d c_i^{(j)} x^i$ for some coefficients $c_i^{(j)}$, for all $x \in I_j$.

We construct a cover for piecewise polynomials.

Lemma 6.24. *There exists a universal constant $c > 0$ such that there is an α -cover of the set of (t, d, k) -piecewise polynomials of size*

$$\binom{k}{t-1} \cdot \left(\frac{tk \cdot e^{cd^{1/2}}}{\alpha} \right)^{(d+1)t}.$$

Proof. We specify an element of the cover by

1. Selecting one of $\binom{k}{t-1}$ partitions of $[k]$ into t intervals I_1, \dots, I_t , and
2. For each interval I_j , selecting an element of an (α/t) -cover \mathcal{C}_j of the set of degree- d polynomials over I_j which are uniformly bounded by 1.

The total size of the cover is $\binom{k}{t-1} \prod_{j=1}^t |\mathcal{C}_j|$. The theorem follows from Proposition 6.25 below, which constructs an (α/t) -cover \mathcal{C}_j of size at most $\left(\frac{tk \cdot e^{cd^{1/2}}}{\alpha} \right)^{d+1}$ for every interval I_j . □

Proposition 6.25. *There exist constants $b, c > 0$ for which the following holds. Let $I \subseteq [k]$ be an interval and let \mathcal{P} be the set of polynomials $p : I \rightarrow \mathbb{R}$ of degree d such that $|p(x)| \leq 1$ for all $x \in I$. There exists an α -cover of \mathcal{P} of size*

$$\min \left\{ \left(\frac{2k}{\alpha} \right)^{|I|}, \left(\frac{ckd^2 \cdot e^{bd^2/|I|}}{\alpha} \right)^{d+1} \right\}.$$

The proof of Proposition 6.25 relies on two major results in approximation theory, which we now state.

Lemma 6.26 (Duffin and Schaeffer [DS41]). *Let $p : [-1, 1] \rightarrow \mathbb{R}$ be a polynomial such that $|p(x)| \leq 1$ for all x of the form $x = \cos(j\pi/d)$ for $j = 0, 1, \dots, d$. Then $|p'(x)| \leq d^2$ for all $x \in [-1, 1]$.*

Lemma 6.27 (Coppersmith and Rivlin [CR92]). *There exist constants $a, b > 0$ for which the following holds. Let $p : \mathbb{R} \rightarrow \mathbb{R}$ be a polynomial of degree d , and suppose that $|p(t)| \leq 1$ for all $t = 0, 1, \dots, m$. Then $|p(t)| \leq a \exp(bd^2/m)$ for all $t \in [0, m]$.*

Proof of Proposition 6.25. We consider two cases, corresponding to the two terms in the minimum. First, consider the function $f : I \rightarrow \mathbb{R}$ where $f(t)$ is obtained by rounding $p(t)$ to the nearest multiple of α/k . Then f satisfies $\sum_{t \in I} |f(t) - p(t)| \leq \alpha$. There are at most $(2k/\alpha)^{|I|}$ functions f which can be constructed this way, giving the first term in the maximum.

For the second term, we construct a cover for \mathcal{P} by approximately interpolating through $d+1$ carefully chosen points in the continuous interval corresponding to I . By applying an affine shift, we may assume that $I = \{0, 1, \dots, m\}$ for some integer $m \leq k-1$. Let $p \in \mathcal{P}$ and for $x \in [0, m]$ let $\hat{p}(x)$ be the value of $p(x)$ rounded to the nearest integer multiple of $\alpha/(2kd^2)$. Let $q : [0, m] \rightarrow \mathbb{R}$

be the unique degree- d polynomial obtained by interpolating through the points $(x_j, \hat{p}(x_j))$ where $x_j = m(1 + \cos(j\pi/d))/2$ for $j = 0, 1, \dots, d$.

We first argue that the polynomial q so defined satisfies $\sum_{t=0}^m |p(t) - q(t)| \leq \alpha$. Let $r(x) = p(x) - q(x)$ for $x \in [0, m]$. Then by construction, $|r(x_j)| \leq \alpha/(2kd^2)$ for all interpolation points x_j . By the Duffin-Schaeffer Inequality (Lemma 6.26), we therefore have $|r'(x)| \leq \frac{\alpha}{km}$ for all $x \in [0, m]$. By the Fundamental Theorem of Calculus, $r(t) = r(0) + \int_0^t r'(t) dt$ satisfies $|r(t)| \leq (t+1) \cdot \frac{\alpha}{km} \leq \alpha/k$, and hence $\sum_{t=0}^m |r(t)| \leq \alpha$.

We now argue that the set of polynomials q that can be constructed in this fashion has size $(ckd^2 \exp(bd^2/m)/\alpha)^{d+1}$. By the Coppersmith-Rivlin Inequality (Lemma 6.27), there are constants $a, b > 0$ such that $|p(x)| \leq a \exp(bd^2/m)$ for all $x \in [0, m]$. Therefore, for each $p \in \mathcal{P}$ and each interpolation point x_j , there are at most $4a \cdot kd^2 \exp(bd^2/m)/\alpha$ possible values that $\hat{p}(x_j)$ can take. Hence, the polynomial q can take one of at most $(4a \cdot kd^2 \exp(bd^2/m)/\alpha)^{d+1}$ possible values, as we wanted to show. \square

Lemma 6.28. *The VC dimension of (t, d, k) -piecewise polynomial distributions is at most $2t(d+1)$.*

Proof. Consider two piecewise polynomial distributions. The difference between their probability mass functions is a piecewise polynomial of degree $\leq d$. The number of intervals needed to represent this piecewise function is $\leq 2t$. It follows that this difference can change sign at most $2td + 2t - 1$ times – each polynomial can change sign at most d times and the sign can change at the interval boundaries. Thus such a function cannot label $2td + 2t + 1$ points with alternating signs, which implies the VC bound. \square

As a corollary, we obtain the following learning algorithm.

Corollary 6.29. *Suppose we are given a set of samples $X_1, \dots, X_n \sim P$, where P is α -close to a (t, d, k) -piecewise polynomial. Then there exists an ε -differentially private algorithm which outputs a (t, d, k) -piecewise polynomial H^* such that $d_{\text{TV}}(P, H^*) \leq (6 + 2\zeta)\alpha$ with probability $\geq 9/10$, so long as*

$$n = \Omega \left(\frac{(d+1)t}{\alpha^2} + \frac{(d+1)t}{\alpha\varepsilon} \cdot \left(\sqrt{d+1} \log k + \log \left(\frac{t}{\alpha} \right) \right) \right).$$

We compare with the work of Diakonikolas, Hardt, and Schmidt [DHS15]. They present an efficient algorithm for $(t, 1, k)$ -piecewise polynomials, with sample complexity $\tilde{O} \left(\frac{t}{\alpha^2} + \frac{t \log k}{\alpha\varepsilon} \right)$, which our algorithm matches⁶. They also claim their results extend to (t, d, k) -piecewise polynomials, though no theorem statement is provided. While we have not investigated the details of this extension, we believe the resulting sample complexity should be qualitatively similar to ours, plausibly with the factor of $\frac{t(d+1)^{3/2} \log k}{\alpha\varepsilon}$ replaced by $\frac{t(d+1) \log k}{\alpha\varepsilon}$.

6.5 Mixtures

In this section, we show that our results immediately extend to learning mixtures of classes of distributions.

Definition 6.30. *Let \mathcal{H} be some set of distributions. A k -mixture of \mathcal{H} is a distribution with density $\sum_{i=1}^k w_i P_i$, where each $P_i \in \mathcal{H}$.*

⁶As stated in [DHS15], their algorithm guarantees approximate differential privacy, but swapping in an appropriate pure DP subroutine gives this result.

Our results follow roughly due to the fact that a cover for k -mixtures of a class can be written as the Cartesian product of k covers for the class. More precisely, we state the following result which bounds the size of the cover of the set of k -mixtures.

Lemma 6.31. *Consider the class of k -mixtures of \mathcal{H} , where \mathcal{H} is some set of distributions. There exists a 2α -cover of this class of size $|\mathcal{C}_\alpha|^k \left(\frac{k}{2\alpha} + 1\right)^{k-1}$, where \mathcal{C}_α is an α -cover of \mathcal{H} .*

Proof. Each element in the cover of the class of mixtures will be obtained by taking k distributions from \mathcal{C}_α , in combination with k mixing weights, which are selected from the set $\{0, \frac{2\alpha}{k}, \frac{4\alpha}{k}, \dots, 1\}$, such that the sum of the mixing weights is 1. The size of this cover is $|\mathcal{C}_\alpha|^k \cdot \left(\frac{k}{2\alpha} + 1\right)^{k-1}$. We reason about the accuracy of the cover as follows. Fix some mixture of k distributions as $\sum_{i=1}^k w_i^{(1)} P_i^{(1)}$, and we will reason about the closest element in our cover, $\sum_{i=1}^k w_i^{(2)} P_i^{(2)}$. By triangle inequality, we have that

$$d_{\text{TV}} \left(\sum_{i=1}^k w_i^{(1)} P_i^{(1)}, \sum_{i=1}^k w_i^{(2)} P_i^{(2)} \right) \leq \sum_{i=1}^k \frac{1}{2} \left| w_i^{(1)} - w_i^{(2)} \right| + w_i^{(1)} d_{\text{TV}} \left(P_i^{(1)}, P_i^{(2)} \right).$$

Since \mathcal{C}_α is an α -cover and $\sum_{i=1}^k w_i^{(1)} = 1$, the total variation distance incurred by the second term will be at most α . As for the mixing weights, note that for the first $k-1$ weights, the nearest weight is at distance at most $\frac{\alpha}{k}$, contributing a total of less than $\frac{\alpha}{2}$. The last mixing weight can be rewritten in terms of the sum of the errors of the other mixing weights, similarly contributing another total of less than $\frac{\alpha}{2}$. This results in the total error being at most 2α , as desired. \square

With this in hand, the following corollary is almost immediate from Corollary 1.2. The factor of $(9 + 3\zeta)\alpha$ (as opposed to $(6 + 2\zeta)\alpha$) is because the closest distribution in the cover of mixture distributions is 3α -close to be P (rather than 2α).

Corollary 6.32. *Let $X_1, \dots, X_n \sim P$, where P is α -close to a k -mixture of distributions from some set \mathcal{H} . Let \mathcal{C}_α be an α -cover of the set \mathcal{H} , and $\zeta > 0$ be a constant. There exists an ε -differentially private algorithm which outputs a distribution which is $(9 + 3\zeta)\alpha$ -close to P with probability $\geq 9/10$, as long as*

$$n = \Omega \left((k \log |\mathcal{C}_\alpha| + k \log(k/\alpha)) \left(\frac{1}{\alpha^2} + \frac{1}{\alpha\varepsilon} \right) \right).$$

For example, instantiating this for mixtures of Gaussians (and disregarding terms which depend on R and κ), we get an algorithm with sample complexity $\tilde{O} \left(\frac{kd^2}{\alpha^2} + \frac{kd^2}{\alpha\varepsilon} \right)$.

6.6 Supervised Learning

We describe an application of our results to the task of binary classification, as modeled by differentially private PAC learning [KLN⁺11]. Let $\mathcal{F} = \{f : X \rightarrow \{0, 1\}\}$ be a publicly known *concept class* of Boolean functions over a domain X . Let P be an unknown probability distribution over X , and let f be an unknown function from \mathcal{F} . Given a sequence $\{(x_i, f(x_i))\}_{i=1}^n$ of i.i.d. samples from P together with their labels under f , the goal of a PAC learner L is to identify a hypothesis $h : X \rightarrow \{0, 1\}$ such that $\Pr_{x \sim P}[h(x) \neq f(x)] \leq \alpha$ for some error parameter $\alpha > 0$. We say that L is (α, β) -accurate if for every $f \in \mathcal{F}$ and every distribution P , it is able to identify such a hypothesis h with probability at least $1 - \beta$ over the choice of the sample and any internal randomness of L .

One of the core results of statistical learning theory is that the sample complexity of *non-private* PAC learning is characterized, up to constant factors, by the VC dimension of the concept class \mathcal{F} .

When one additionally requires the learner L to be differentially private with respect to its input sample, such a characterization is unknown. However, it is known that the sample complexity of private learning can be arbitrarily higher than that of non-private learning. For example, when $\mathcal{F} = \{f_t : t \in X\}$ is the class of threshold functions defined by $f_t(x) = 1 \iff x \leq t$ over a totally ordered domain X , the sample complexity of PAC learning under the most permissive notion of (ε, δ) -differential privacy is $\Omega(\log^* |X|)$ [BNSV15, ALMM19]. Meanwhile, the VC dimension of this class, and hence the sample complexity of non-private learning, is a constant independent of $|X|$.

While this separation shows that there can be a sample cost of privacy for PAC learning, this cost can be completely eliminated if the distribution P on examples is known. This was observed by Beimel, Nissim, and Stemmer [BNS16], who showed that if a good approximation to P is known, e.g., from public unlabeled examples or from differentially private processing of unlabeled examples, then the number of labeled examples needed for private PAC learning is only $O(VC(\mathcal{F}))$.

Theorem 6.33. *Let $\varepsilon > 0$, $\mathcal{F} = \{f : X \rightarrow \{0, 1\}\}$, and P be a publicly known distribution over X . For $n = O\left(\frac{1}{\alpha^2 \varepsilon}(VC(\mathcal{F}) \log(1/\alpha) + \log(1/\beta))\right)$, there exists an ε -differentially private algorithm $L : (X \times \{0, 1\})^n \rightarrow \mathcal{F}$ such that for every $f \in \mathcal{F}$, with probability at least $1 - \beta$ over the choice of $x_1, \dots, x_n \leftarrow P$, we have that $L((x_1, f(x_1)), \dots, (x_n, f(x_n)))$ produces $h \in \mathcal{F}$ such that $\Pr_{x \sim P}[f(x) \neq h(x)] \leq \alpha$.*

Our results suggest a natural two-step algorithm for private PAC learning when the distribution P itself is not known, but is known to (approximately) come from a set of distributions \mathcal{H} : The algorithm first uses private hypothesis selection to select \hat{H} with $d_{TV}(P, \hat{H}) \leq \alpha/2$, and then runs the algorithm of [BNS16] using \hat{H} in place of P with error parameter $\alpha/2$. Using the fact that $d_{TV}(P, \hat{H}) \leq \alpha/2$ implies $|\Pr_{x \sim P}[f(x) \neq h(x)] - \Pr_{x \sim \hat{H}}[f(x) \neq h(x)]| \leq \alpha/2$, the following result holds by combining Theorem 6.33 with Corollary 1.2.

Corollary 6.34. *Let \mathcal{H} be a set of distributions over X with an α -cover \mathcal{C}_α . Let P be a distribution over X with $d_{TV}(P, \mathcal{H}) \leq \alpha/(4(3 + \zeta))$. Then for*

$$n = O\left(\frac{\log |\mathcal{C}_\alpha|}{\alpha^2} + \frac{\log |\mathcal{C}_\alpha|}{\alpha \varepsilon} + \frac{VC(\mathcal{F}) \log(1/\alpha)}{\alpha^2 \varepsilon}\right)$$

there exists an ε -differentially private algorithm $L : (X \times \{0, 1\})^n \rightarrow \mathcal{F}$ such that for every $f \in \mathcal{F}$, with probability at least $3/4$ over the choice of $x_1, \dots, x_n \leftarrow P$, we have that $L((x_1, f(x_1)), \dots, (x_n, f(x_n)))$ produces $h \in \mathcal{F}$ such that $\Pr_{x \sim P}[f(x) \neq h(x)] \leq \alpha$.

Theorem 6.33 can, of course, also be combined with the more refined guarantees of Theorem 4.1. As an example application, combining Theorem 6.33 with Corollary 6.13 gives a (ε, δ) -differentially private algorithm for learning one-dimension thresholds with respect to univariate Gaussian distributions on the reals. In contrast, this task is impossible without making distributional assumptions.

7 Conclusions

In this paper, we presented differentially private methods for hypothesis selection. The sample complexity can be bounded by the logarithm of the number of hypotheses. This allows us to provide bounds on the sample complexity of (semi-agnostically) learning a class which depend on the logarithm of the covering number, complementing known lower bounds which depend on the logarithm of the packing number. There are many interesting questions left open by our work, a few of which we outline below.

1. Our algorithms for learning classes of distributions all use cover-based arguments, and thus are not computationally efficient. For instance, we provide the first $\tilde{O}(d)$ sample complexity upper bound on ε -differentially privately learning a product distribution and Gaussian with known covariance. One interesting question is whether there is an efficient algorithm which achieves this sample complexity.
2. The running time of our method is quadratic in the number of hypotheses – is it possible to reduce this to a near-linear time complexity?
3. Our main theorem obtains an approximation factor which is arbitrarily close to 3, which is optimal for this problem, even without privacy. This factor can be reduced to 2 if one is OK with outputting a mixture of hypotheses from the set [BKM19]. Is this achievable with privacy constraints?

Acknowledgments

The authors would like to thank Shay Moran for bringing to their attention the application to PAC learning mentioned in Section 6.6, Jonathan Ullman for asking questions which motivated Remark 6.4, and Clément Canonne for assistance in reducing the constant factor in the approximation guarantee.

References

- [ABDH⁺18] Hassan Ashtiani, Shai Ben-David, Nicholas Harvey, Christopher Liaw, Abbas Mehrabian, and Yaniv Plan. Nearly tight sample complexity bounds for learning mixtures of Gaussians via sample compression schemes. In *Advances in Neural Information Processing Systems 31*, NeurIPS '18, pages 3412–3421. Curran Associates, Inc., 2018.
- [ABG⁺14] Joseph Anderson, Mikhail Belkin, Navin Goyal, Luis Rademacher, and James R. Voss. The more, the merrier: the blessing of dimensionality for learning large Gaussian mixtures. In *Proceedings of the 27th Annual Conference on Learning Theory*, COLT '14, pages 1135–1164, 2014.
- [ADK15] Jayadev Acharya, Constantinos Daskalakis, and Gautam Kamath. Optimal testing for properties of distributions. In *Advances in Neural Information Processing Systems 28*, NIPS '15, pages 3577–3598. Curran Associates, Inc., 2015.
- [ADLS17] Jayadev Acharya, Ilias Diakonikolas, Jerry Li, and Ludwig Schmidt. Sample-optimal density estimation in nearly-linear time. In *Proceedings of the 28th Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '17, pages 1278–1289, Philadelphia, PA, USA, 2017. SIAM.
- [AFJ⁺18] Jayadev Acharya, Moein Falahatgar, Ashkan Jafarpour, Alon Orlitsky, and Ananda Theertha Suresh. Maximum selection and sorting with adversarial comparators. *Journal of Machine Learning Research*, 19(1):2427–2457, 2018.
- [AJOS14] Jayadev Acharya, Ashkan Jafarpour, Alon Orlitsky, and Ananda Theertha Suresh. Sorting with adversarial comparators and application to density estimation. In *Proceedings of the 2014 IEEE International Symposium on Information Theory*, ISIT '14, pages 1682–1686, Washington, DC, USA, 2014. IEEE Computer Society.

- [AK01] Sanjeev Arora and Ravi Kannan. Learning mixtures of arbitrary Gaussians. In *Proceedings of the 33rd Annual ACM Symposium on the Theory of Computing*, STOC '01, pages 247–257, New York, NY, USA, 2001. ACM.
- [AKSZ18] Jayadev Acharya, Gautam Kamath, Ziteng Sun, and Huanyu Zhang. Inspectre: Privately estimating the unseen. In *Proceedings of the 35th International Conference on Machine Learning*, ICML '18, pages 30–39. JMLR, Inc., 2018.
- [ALMM19] Noga Alon, Roi Livni, Maryanthe Malliaris, and Shay Moran. Private PAC learning implies finite Littlestone dimension. In *Proceedings of the 51st Annual ACM Symposium on the Theory of Computing*, STOC '19, New York, NY, USA, 2019. ACM.
- [AM05] Dimitris Achlioptas and Frank McSherry. On spectral learning of mixtures of distributions. In *Proceedings of the 18th Annual Conference on Learning Theory*, COLT '05, pages 458–469. Springer, 2005.
- [Ant95] Martin Anthony. Classification by polynomial surfaces. *Discrete Applied Mathematics*, 61(2):91–103, 1995.
- [AS12] Pranjal Awasthi and Or Sheffet. Improved spectral-norm bounds for clustering. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques.*, APPROX '12, pages 37–49. Springer, 2012.
- [ASZ19] Jayadev Acharya, Ziteng Sun, and Huanyu Zhang. Hadamard response: Estimating distributions privately, efficiently, and with little communication. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, AISTATS '19, pages 1120–1129. JMLR, Inc., 2019.
- [BBKN14] Amos Beimel, Hai Brenner, Shiva Prasad Kasiviswanathan, and Kobbi Nissim. Bounds on the sample complexity for private learning and private data release. *Machine Learning*, 94(3):401–437, 2014.
- [BCM14] Aditya Bhaskara, Moses Charikar, Ankur Moitra, and Aravindan Vijayaraghavan. Smoothed analysis of tensor decompositions. In *Proceedings of the 46th Annual ACM Symposium on the Theory of Computing*, STOC '14, pages 594–603, New York, NY, USA, 2014. ACM.
- [BDMN05] Avrim Blum, Cynthia Dwork, Frank McSherry, and Kobbi Nissim. Practical privacy: The SuLQ framework. In *Proceedings of the 24th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, PODS '05, pages 128–138, New York, NY, USA, 2005. ACM.
- [BDRS18] Mark Bun, Cynthia Dwork, Guy N. Rothblum, and Thomas Steinke. Composable and versatile privacy via truncated cdp. In *Proceedings of the 50th Annual ACM Symposium on the Theory of Computing*, STOC '18, pages 74–86, New York, NY, USA, 2018. ACM.
- [BKM19] Olivier Bousquet, Daniel M. Kane, and Shay Moran. The optimal approximation factor in density estimation. In *Proceedings of the 32nd Annual Conference on Learning Theory*, COLT '19, pages 318–341, 2019.
- [BNS16] Amos Beimel, Kobbi Nissim, and Uri Stemmer. Private learning and sanitization: Pure vs. approximate differential privacy. *Theory of Computing*, 12(1):1–61, 2016.

- [BNSV15] Mark Bun, Kobbi Nissim, Uri Stemmer, and Salil Vadhan. Differentially private release and learning of threshold functions. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '15, pages 634–649, Washington, DC, USA, 2015. IEEE Computer Society.
- [BS10] Mikhail Belkin and Kaushik Sinha. Polynomial learning of distribution families. In *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science*, FOCS '10, pages 103–112, Washington, DC, USA, 2010. IEEE Computer Society.
- [BS16] Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Proceedings of the 14th Conference on Theory of Cryptography*, TCC '16-B, pages 635–658, Berlin, Heidelberg, 2016. Springer.
- [BUV14] Mark Bun, Jonathan Ullman, and Salil Vadhan. Fingerprinting codes and the price of approximate differential privacy. In *Proceedings of the 46th Annual ACM Symposium on the Theory of Computing*, STOC '14, pages 1–10, New York, NY, USA, 2014. ACM.
- [CDSS14a] Siu On Chan, Ilias Diakonikolas, Rocco A. Servedio, and Xiaorui Sun. Efficient density estimation via piecewise polynomial approximation. In *Proceedings of the 46th Annual ACM Symposium on the Theory of Computing*, STOC '14, pages 604–613, New York, NY, USA, 2014. ACM.
- [CDSS14b] Siu On Chan, Ilias Diakonikolas, Rocco A. Servedio, and Xiaorui Sun. Near-optimal density estimation in near-linear time using variable-width histograms. In *Advances in Neural Information Processing Systems 27*, NIPS '14, pages 1844–1852. Curran Associates, Inc., 2014.
- [CR92] Don Coppersmith and Theodore J. Rivlin. The growth of polynomials bounded at equally spaced points. *SIAM Journal on Mathematical Analysis*, 23(4):970–983, 1992.
- [CR08a] Kamalika Chaudhuri and Satish Rao. Beyond Gaussians: Spectral methods for learning mixtures of heavy-tailed distributions. In *Proceedings of the 21st Annual Conference on Learning Theory*, COLT '08, pages 21–32, 2008.
- [CR08b] Kamalika Chaudhuri and Satish Rao. Learning mixtures of product distributions using correlations and independence. In *Proceedings of the 21st Annual Conference on Learning Theory*, COLT '08, pages 9–20, 2008.
- [CWZ19] T. Tony Cai, Yichen Wang, and Linjun Zhang. The cost of privacy: Optimal rates of convergence for parameter estimation with differential privacy. *arXiv preprint arXiv:1902.04495*, 2019.
- [Das99] Sanjoy Dasgupta. Learning mixtures of Gaussians. In *Proceedings of the 40th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '99, pages 634–644, Washington, DC, USA, 1999. IEEE Computer Society.
- [DDKT16] Constantinos Daskalakis, Anindya De, Gautam Kamath, and Christos Tzamos. A size-free CLT for Poisson multinomials and its applications. In *Proceedings of the 48th Annual ACM Symposium on the Theory of Computing*, STOC '16, pages 1074–1086, New York, NY, USA, 2016. ACM.

- [DDO⁺13] Constantinos Daskalakis, Ilias Diakonikolas, Ryan O'Donnell, Rocco A. Servedio, and Li Yang Tan. Learning sums of independent integer random variables. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '13, pages 217–226, Washington, DC, USA, 2013. IEEE Computer Society.
- [DDS12a] Constantinos Daskalakis, Ilias Diakonikolas, and Rocco A. Servedio. Learning k-modal distributions via testing. In *Proceedings of the 23th Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '12, pages 1371–1385, Philadelphia, PA, USA, 2012. SIAM.
- [DDS12b] Constantinos Daskalakis, Ilias Diakonikolas, and Rocco A. Servedio. Learning Poisson binomial distributions. In *Proceedings of the 44th Annual ACM Symposium on the Theory of Computing*, STOC '12, pages 709–728, New York, NY, USA, 2012. ACM.
- [DHS15] Ilias Diakonikolas, Moritz Hardt, and Ludwig Schmidt. Differentially private learning of structured discrete distributions. In *Advances in Neural Information Processing Systems 28*, NIPS '15, pages 2566–2574. Curran Associates, Inc., 2015.
- [Dif17] Differential Privacy Team, Apple. Learning with privacy at scale. <https://machinelearning.apple.com/docs/learning-with-privacy-at-scale/appliedifferentialprivacysystem.pdf>, December 2017.
- [DJW13] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. Local privacy and statistical minimax rates. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '13, pages 429–438, Washington, DC, USA, 2013. IEEE Computer Society.
- [DK14] Constantinos Daskalakis and Gautam Kamath. Faster and sample near-optimal algorithms for proper learning mixtures of Gaussians. In *Proceedings of the 27th Annual Conference on Learning Theory*, COLT '14, pages 1183–1213, 2014.
- [DKK⁺16] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high dimensions without the computational intractability. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '16, pages 655–664, Washington, DC, USA, 2016. IEEE Computer Society.
- [DKS16a] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. The Fourier transform of Poisson multinomial distributions and its algorithmic applications. In *Proceedings of the 48th Annual ACM Symposium on the Theory of Computing*, STOC '16, New York, NY, USA, 2016. ACM.
- [DKS16b] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. Optimal learning via the Fourier transform for sums of independent integer random variables. In *Proceedings of the 29th Annual Conference on Learning Theory*, COLT '16, pages 831–849, 2016.
- [DKS16c] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. Properly learning Poisson binomial distributions in almost polynomial time. In *Proceedings of the 29th Annual Conference on Learning Theory*, COLT '16, pages 850–878, 2016.
- [DKS18] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. List-decodable robust mean estimation and learning mixtures of spherical Gaussians. In *Proceedings of the 50th*

Annual ACM Symposium on the Theory of Computing, STOC '18, pages 1047–1060, New York, NY, USA, 2018. ACM.

- [DKT15] Constantinos Daskalakis, Gautam Kamath, and Christos Tzamos. On the structure, covering, and learning of Poisson multinomial distributions. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '15, pages 1203–1217, Washington, DC, USA, 2015. IEEE Computer Society.
- [DL96] Luc Devroye and Gábor Lugosi. A universally acceptable smoothing factor for kernel density estimation. *The Annals of Statistics*, 24(6):2499–2512, 1996.
- [DL97] Luc Devroye and Gábor Lugosi. Nonasymptotic universal smoothing factors, kernel complexity and Yatracos classes. *The Annals of Statistics*, 25(6):2626–2637, 1997.
- [DL01] Luc Devroye and Gábor Lugosi. *Combinatorial methods in density estimation*. Springer, 2001.
- [DL09] Cynthia Dwork and Jing Lei. Differential privacy and robust statistics. In *Proceedings of the 41st Annual ACM Symposium on the Theory of Computing*, STOC '09, pages 371–380, New York, NY, USA, 2009. ACM.
- [DLS⁺17] Aref N. Dajani, Amy D. Lauger, Phyllis E. Singer, Daniel Kifer, Jerome P. Reiter, Ashwin Machanavajjhala, Simson L. Garfinkel, Scot A. Dahl, Matthew Graham, Vishesh Karwa, Hang Kim, Philip Lelerc, Ian M. Schmutte, William N. Sexton, Lars Vilhuber, and John M. Abowd. The modernization of statistical disclosure limitation at the U.S. census bureau, 2017. Presented at the September 2017 meeting of the Census Scientific Advisory Committee.
- [DLS18] Anindya De, Philip M. Long, and Rocco A. Servedio. Learning sums of independent random variables with sparse collective support. In *Proceedings of the 59th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '18, pages 297–308, Washington, DC, USA, 2018. IEEE Computer Society.
- [DMNS06] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the 3rd Conference on Theory of Cryptography*, TCC '06, pages 265–284, Berlin, Heidelberg, 2006. Springer.
- [DMR18] Luc Devroye, Abbas Mehrabian, and Tommy Reddad. The total variation distance between high-dimensional Gaussians. *arXiv preprint arXiv:1810.08693*, 2018.
- [DP08] Constantinos Daskalakis and Christos H. Papadimitriou. Discretized multinomial distributions and Nash equilibria in anonymous games. In *Proceedings of the 49th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '08, pages 25–34, Washington, DC, USA, 2008. IEEE Computer Society.
- [DP09] Constantinos Daskalakis and Christos H. Papadimitriou. On oblivious PTAS's for Nash equilibrium. In *Proceedings of the 41st Annual ACM Symposium on the Theory of Computing*, STOC '09, pages 75–84, New York, NY, USA, 2009. ACM.
- [DP15a] Constantinos Daskalakis and Christos H. Papadimitriou. Approximate Nash equilibria in anonymous games. *Journal of Economic Theory*, 156:207–245, 2015.

- [DP15b] Constantinos Daskalakis and Christos H. Papadimitriou. Sparse covers for sums of indicators. *Probability Theory and Related Fields*, 162(3):679–705, 2015.
- [DR16] Cynthia Dwork and Guy N. Rothblum. Concentrated differential privacy. *arXiv preprint arXiv:1603.01887*, 2016.
- [DR18] John C. Duchi and Feng Ruan. The right complexity measure in locally private estimation: It is not the fisher information. *arXiv preprint arXiv:1806.05756*, 2018.
- [DS41] Richard J. Duffin and Albert C. Schaeffer. A refinement of an inequality of the brothers Markoff. *Transactions of the American Mathematical Society*, 50(3):517–528, 1941.
- [DS00] Sanjoy Dasgupta and Leonard J. Schulman. A two-round variant of EM for Gaussian mixtures. In *Proceedings of the 16th Conference in Uncertainty in Artificial Intelligence*, UAI '00, pages 152–159. Morgan Kaufmann, 2000.
- [DTZ17] Constantinos Daskalakis, Christos Tzamos, and Manolis Zampetakis. Ten steps of EM suffice for mixtures of two Gaussians. In *Proceedings of the 30th Annual Conference on Learning Theory*, COLT '17, pages 704–710, 2017.
- [EPK14] Úlfar Erlingsson, Vasyl Pihur, and Aleksandra Korolova. RAPPOR: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 ACM Conference on Computer and Communications Security*, CCS '14, pages 1054–1067, New York, NY, USA, 2014. ACM.
- [FOS06] Jon Feldman, Ryan O'Donnell, and Rocco A. Servedio. PAC learning axis-aligned mixtures of Gaussians with no separation assumption. In *Proceedings of the 19th Annual Conference on Learning Theory*, COLT '06, pages 20–34, Berlin, Heidelberg, 2006. Springer.
- [FOS08] Jon Feldman, Ryan O'Donnell, and Rocco A. Servedio. Learning mixtures of product distributions over discrete domains. *SIAM Journal on Computing*, 37(5):1536–1564, 2008.
- [GHK15] Rong Ge, Qingqing Huang, and Sham M. Kakade. Learning mixtures of Gaussians in high dimensions. In *Proceedings of the 47th Annual ACM Symposium on the Theory of Computing*, STOC '15, pages 761–770, New York, NY, USA, 2015. ACM.
- [GRS19] Marco Gaboardi, Ryan Rogers, and Or Sheffet. Locally private confidence intervals: Z-test and tight confidence intervals. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, AISTATS '19, pages 2545–2554. JMLR, Inc., 2019.
- [HK13] Daniel Hsu and Sham M. Kakade. Learning mixtures of spherical Gaussians: Moment methods and spectral decompositions. In *Proceedings of the 4th Conference on Innovations in Theoretical Computer Science*, ITCS '13, pages 11–20, New York, NY, USA, 2013. ACM.
- [HL18] Samuel B. Hopkins and Jerry Li. Mixture models, robustness, and sum of squares proofs. In *Proceedings of the 50th Annual ACM Symposium on the Theory of Computing*, STOC '18, pages 1021–1034, New York, NY, USA, 2018. ACM.

- [HP15] Moritz Hardt and Eric Price. Sharp bounds for learning a mixture of two Gaussians. In *Proceedings of the 47th Annual ACM Symposium on the Theory of Computing*, STOC '15, pages 753–760, New York, NY, USA, 2015. ACM.
- [HT10] Moritz Hardt and Kunal Talwar. On the geometry of differential privacy. In *Proceedings of the 42nd Annual ACM Symposium on the Theory of Computing*, STOC '10, pages 705–714, New York, NY, USA, 2010. ACM.
- [JKMW18] Matthew Joseph, Janardhan Kulkarni, Jieming Mao, and Zhiwei Steven Wu. Locally private Gaussian estimation. *arXiv preprint arXiv:1811.08382*, 2018.
- [KBR16] Peter Kairouz, Keith Bonawitz, and Daniel Ramage. Discrete distribution estimation under local privacy. In *Proceedings of the 33rd International Conference on Machine Learning*, ICML '16, pages 2436–2444. JMLR, Inc., 2016.
- [KK10] Amit Kumar and Ravindran Kannan. Clustering with spectral norm and the k-means algorithm. In *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science*, FOCS '10, pages 299–308, Washington, DC, USA, 2010. IEEE Computer Society.
- [KLN⁺11] Shiva Prasad Kasiviswanathan, Homin K. Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. What can we learn privately? *SIAM Journal on Computing*, 40(3):793–826, 2011.
- [KLSU19] Gautam Kamath, Jerry Li, Vikrant Singhal, and Jonathan Ullman. Privately learning high-dimensional distributions. In *Proceedings of the 32nd Annual Conference on Learning Theory*, COLT '19, pages 1853–1902, 2019.
- [KMR⁺94] Michael Kearns, Yishay Mansour, Dana Ron, Ronitt Rubinfeld, Robert E. Schapire, and Linda Sellie. On the learnability of discrete distributions. In *Proceedings of the 26th Annual ACM Symposium on the Theory of Computing*, STOC '94, pages 273–282, New York, NY, USA, 1994. ACM.
- [KMV10] Adam Tauman Kalai, Ankur Moitra, and Gregory Valiant. Efficiently learning mixtures of two Gaussians. In *Proceedings of the 42nd Annual ACM Symposium on the Theory of Computing*, STOC '10, pages 553–562, New York, NY, USA, 2010. ACM.
- [KSS18] Pravesh Kothari, Jacob Steinhardt, and David Steurer. Robust moment estimation and improved clustering via sum of squares. In *Proceedings of the 50th Annual ACM Symposium on the Theory of Computing*, STOC '18, pages 1035–1046, New York, NY, USA, 2018. ACM.
- [KSSU19] Gautam Kamath, Or Sheffet, Vikrant Singhal, and Jonathan Ullman. Differentially private algorithms for learning mixtures of separated gaussians. In *Advances in Neural Information Processing Systems 32*, NeurIPS '19. Curran Associates, Inc., 2019.
- [KV18] Vishesh Karwa and Salil Vadhan. Finite sample differentially private confidence intervals. In *Proceedings of the 9th Conference on Innovations in Theoretical Computer Science*, ITCS '18, pages 44:1–44:9, Dagstuhl, Germany, 2018. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.

- [LS17] Jerry Li and Ludwig Schmidt. Robust proper learning for mixtures of Gaussians via systems of polynomial inequalities. In *Proceedings of the 30th Annual Conference on Learning Theory*, COLT '17, pages 1302–1382, 2017.
- [Mir17] Ilya Mironov. Rényi differential privacy. In *Proceedings of the 30th IEEE Computer Security Foundations Symposium*, CSF '17, pages 263–275, Washington, DC, USA, 2017. IEEE Computer Society.
- [MS08] Satyaki Mahalanabis and Daniel Stefankovic. Density estimation in linear time. In *Proceedings of the 21st Annual Conference on Learning Theory*, COLT '08, pages 503–512, 2008.
- [MT07] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '07, pages 94–103, Washington, DC, USA, 2007. IEEE Computer Society.
- [MV10] Ankur Moitra and Gregory Valiant. Settling the polynomial learnability of mixtures of Gaussians. In *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science*, FOCS '10, pages 93–102, Washington, DC, USA, 2010. IEEE Computer Society.
- [NRS07] Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. Smooth sensitivity and sampling in private data analysis. In *Proceedings of the 39th Annual ACM Symposium on the Theory of Computing*, STOC '07, pages 75–84, New York, NY, USA, 2007. ACM.
- [RV17] Oded Regev and Aravindan Vijayaraghavan. On learning mixtures of well-separated Gaussians. In *Proceedings of the 58th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '17, pages 85–96, Washington, DC, USA, 2017. IEEE Computer Society.
- [Smi11] Adam Smith. Privacy-preserving statistical estimation with optimal convergence rates. In *Proceedings of the 43rd Annual ACM Symposium on the Theory of Computing*, STOC '11, pages 813–822, New York, NY, USA, 2011. ACM.
- [SOAJ14] Ananda Theertha Suresh, Alon Orlitsky, Jayadev Acharya, and Ashkan Jafarpour. Near-optimal-sample estimators for spherical Gaussian mixtures. In *Advances in Neural Information Processing Systems 27*, NIPS '14, pages 1395–1403. Curran Associates, Inc., 2014.
- [SU15] Thomas Steinke and Jonathan Ullman. Interactive fingerprinting codes and the hardness of preventing false discovery. In *Proceedings of the 28th Annual Conference on Learning Theory*, COLT '15, pages 1588–1628, 2015.
- [SU17] Thomas Steinke and Jonathan Ullman. Between pure and approximate differential privacy. *The Journal of Privacy and Confidentiality*, 7(2):3–22, 2017.
- [Tal94] Michel Talagrand. Sharper bounds for Gaussian and empirical processes. *The Annals of Probability*, 22(1):28–76, 1994.
- [TS13] Abhradeep Guha Thakurta and Adam Smith. Differentially private feature selection via stability arguments, and the robustness of the lasso. In *Proceedings of the 26th Annual Conference on Learning Theory*, COLT '13, pages 819–850, 2013.

- [VC74] Vladimir Vapnik and Alexey Chervonenkis. *Theory of Pattern Recognition*. Nauka, 1974.
- [VV10] Gregory Valiant and Paul Valiant. A CLT and tight lower bounds for estimating entropy. *Electronic Colloquium on Computational Complexity (ECCC)*, 17(179), 2010.
- [VW02] Santosh Vempala and Grant Wang. A spectral algorithm for learning mixtures of distributions. In *Proceedings of the 43rd Annual IEEE Symposium on Foundations of Computer Science*, FOCS '02, pages 113–123, Washington, DC, USA, 2002. IEEE Computer Society.
- [WHW⁺16] Shaowei Wang, Liusheng Huang, Pengzhan Wang, Yiwen Nie, Hongli Xu, Wei Yang, Xiang-Yang Li, and Chunming Qiao. Mutual information optimally local private discrete distribution estimation. *arXiv preprint arXiv:1607.08025*, 2016.
- [XHM16] Ji Xu, Daniel J. Hsu, and Arian Maleki. Global analysis of expectation maximization for mixtures of two Gaussians. In *Advances in Neural Information Processing Systems 29*, NIPS '16, pages 2676–2684. Curran Associates, Inc., 2016.
- [Yat85] Yannis G. Yatracos. Rates of convergence of minimum distance estimators and Kolmogorov’s entropy. *The Annals of Statistics*, 13(2):768–774, 1985.
- [YB18] Min Ye and Alexander Barg. Optimal schemes for discrete distribution estimation under locally differential privacy. *IEEE Transactions on Information Theory*, 64(8):5662–5676, 2018.