Implicit High-Order Moment Tensor Estimation and Learning Latent Variable Models

Ilias Diakonikolas* University of Wisconsin-Madison ilias@cs.wisc.edu Daniel M. Kane[†] University of California, San Diego dakane@cs.ucsd.edu

Abstract

We study the task of learning latent-variable models on \mathbb{R}^d with k hidden parameters. A common obstacle towards designing computationally efficient algorithms for such models is the necessity of approximating moment tensors of *super-constant* degree, scaling with the parameter k or the desired accuracy $\epsilon > 0$. Motivated by such learning applications, we develop a general efficient algorithm for *implicit moment tensor computation*. Roughly speaking, our algorithm computes in poly(d, k) time a succinct approximate description of tensors of the form $M_m = \sum_{i=1}^k w_i v_i^{\otimes m}$, for $w_i \in \mathbb{R}_+$ —even for $m = \omega(1)$ —assuming that there exists a *polynomial-size arithmetic circuit* whose expected output on an appropriate samplable distribution is equal to M_m , and whose covariance on this input is bounded. Our framework broadly generalizes, both conceptually and technically, the work of [LL22] which developed an efficient algorithm for the specific moment tensors that arise in the task of clustering mixtures of spherical Gaussians.

By leveraging our general implicit moment estimation algorithm, we obtain the first poly(d, k)time learning algorithms for the following classical latent-variable models—thereby resolving or making significant progress towards a number of important open problems in the literature.

- Mixtures of Linear Regressions Given i.i.d. samples (x, y) with $x \sim N(0, I)$ and such that the joint distribution on (x, y) is an unknown k-mixture of linear regressions on \mathbb{R}^{d+1} corrupted with Gaussian noise, the goal is to learn the underlying distribution in total variation distance. We give a poly $(d, k, 1/\epsilon)$ -time algorithm for this task, where ϵ is the desired error. The previously best algorithm has super-polynomial complexity in k.
- Learning Mixtures of Spherical Gaussians Given i.i.d. samples from a k-mixture of identity covariance Gaussians on \mathbb{R}^d , the goal is to learn the target mixture. We give a poly $(d, k, 1/\epsilon)$ -time density estimation algorithm, where ϵ is the desired total variation error, under the condition that the means lie in a ball of radius $O(\sqrt{\log k})$. Prior algorithms incur super-polynomial complexity in k. For parameter estimation, we give a poly $(d, k, 1/\epsilon)$ -time algorithm where ϵ is the target accuracy, under the *optimal* mean separation of $\Omega(\log^{1/2}(k/\epsilon))$ and the condition that the largest distance is comparable to the smallest. Prior polynomial-time parameter estimation algorithms require separation $\Omega(\log^{1/2+c}(k/\epsilon))$, for c > 0.
- PAC Learning Positive Linear Combinations of ReLUs Given labeled examples (x, y) with $x \sim N(0, I)$ and y = F(x), where F is an unknown positive linear combination of k ReLUs (or other reasonable activations) on \mathbb{R}^d , the goal is to PAC learn the target function in L_2 -norm. We give an algorithm for this task with complexity $poly(d, k)2^{poly(1/\epsilon)}$, where ϵ is the desired relative L_2 -error. Our is the first algorithm for this learning task that runs in fully-polynomial (i.e., poly(d, k)) time for subconstant values of the error $\epsilon = o_{k,d}(1)$. All previously known algorithms have super-polynomial complexity, for any $\epsilon = o_{k,d}(1)$.

^{*}Supported by NSF Medium Award CCF-2107079, and an H.I. Romnes Faculty Fellowship.

[†]Supported by NSF Medium Award CCF-2107547 and NSF Award CCF-1553288 (CAREER).

1 Introduction

This work is motivated by the general algorithmic problem of learning probabilistic models with latent variables in high dimensions. This is a prototypical statistical estimation task, first studied in the pioneering work of Karl Pearson from 1894 [Pea94]. Pearson introduced the method of moments motivated by the problem of learning Gaussian mixtures in one dimension—one of the most basic latent-variable models. The algorithmic version of learning high-dimensional latent variable models, including mixtures of Gaussians, has been extensively studied in recent decades within the theoretical computer science and machine learning communities. The relevant literature is vast and has resulted in significant algorithmic advances for a diverse range of settings.

Here we focus on the regime where both the dimension d and the number of hidden parameters k are large. In this context, we are interested in designing learning algorithms for natural models with complexity poly(d, k), i.e., polynomial in both d and k. Of course, there are examples where such a complexity guarantee may not be possible, due to fundamental computational limitations. One such example is the task of learning (either proper learning or density estimation) k-mixtures of Gaussians with non-spherical covariances on \mathbb{R}^d , for which recent work [DKS17, BRST21, GVV22, DKPP24] has provided formal evidence of computational hardness—ruling out poly(d, k) runtime, even though poly(d, k) samples suffice. For the latent-variable models we study in this paper—namely k-mixtures of linear regressions, k-mixtures of spherical Gaussians, and one-hidden-layer neural networks with k gates— there is no known formal evidence of hardness. On the other hand, the best known algorithms prior to this work had super-polynomial complexity (as a function of k). As a corollary of our main result (Proposition 1.7), we design the first poly(d, k) time learning algorithms for all of these problems. We thus resolve, or make significant progress towards, several well-known open problems in the learning theory literature. All of our algorithmic results rely on a novel methodology that is of broader interest and may find other applications.

Before we describe our main technique, we start by highlighting our algorithmic results on learning latent-variable models.

1.1 Efficiently Learning Latent Variable Models

Mixtures of Linear Regressions A k-mixture of linear regressions (k-MLR), specified by mixing weights $w_i \ge 0$, where $\sum_{i=1}^k w_i = 1$, and regressors $\beta_i \in \mathbb{R}^d$, $i \in [k]$, is the distribution F on pairs $(x, y) \in \mathbb{R}^d \times \mathbb{R}$, where $x \sim N(0, I)$ and with probability w_i we have that $y = \beta_i \cdot x + \nu$, where $\nu \sim N(0, \sigma^2)$ is independent of x. Mixtures of linear regressions are a classical probabilistic model introduced in [DeV89, JJ94] and have since been extensively studied in machine learning; see Section 1.4 for a detailed summary of prior work.

Here we focus on density estimation for k-MLRs, where the goal is to learn the underlying distribution in total variation distance—without making any separation assumptions on the pairwise distance between the components. While density estimation can be information-theoretically solved with poly(d, k) samples, all prior algorithms required super-polynomial time. Specifically, [LL18] gave an algorithm with complexity exponential in k. This bound was improved by [CLS20] to sub-exponential, namely scaling with $\exp(\tilde{O}(k^{1/2}))$. The fastest previously known algorithm [DK20] has quasi-polynomial complexity in k, namely poly $(d, (k/\epsilon)^{\log^2(k)})$, where ϵ is the accuracy.

Here we give the first polynomial-time algorithm. Specifically, we show (see Theorem 4.19):

Theorem 1.1 (Density Estimation for k-MLR). Let F be a k-MLR distribution on \mathbb{R}^{d+1} with $B, \sigma \leq 1$, where $B = \max_i \|\beta_i\|_2$ and σ is the standard deviation of the Gaussian noise. There exists an algorithm that draws $N = \text{poly}(k, d)(1/\epsilon)^{O(\sigma^{-2})}$ samples from F, runs in poly(N, d) time, and returns a sampler for a distribution that is ϵ -close to F in total variation distance.

The existence of a learning algorithm for k-mixtures of linear regressions with a poly(d, k) complexity was a well-known open problem in the literature. This question was most recently posed as one of the main open problems in [CLS20]; see the associated talk at STOC 2020 [Che20]. Theorem 1.1 answers this open question in the affirmative.

It should also be noted that the requirement that $\sigma, B \leq 1$ can be removed at the cost of making the runtime $poly(k, d)(1/\epsilon)^{1+O((B/\sigma)^2)}$. This is done by computing an upper bound A on $B + \sigma$ (for example by taking many samples and returning something proportional to the largest value of |y|), and renormalizing by replacing y by y/A. An interesting open question that remains is to remove the dependence on σ, B from the complexity of our algorithm.

Mixtures of Spherical Gaussians A k-mixture of spherical Gaussians is any distribution on \mathbb{R}^d of the form $F = \sum_{i=1}^k w_i N(\mu_i, I)$, where $\mu_i \in \mathbb{R}^d$ are the unknown mean vectors and $w_i \ge 0$, with $\sum_{i=1}^k w_i = 1$, are the mixing weights. Gaussian mixture models are one of the oldest and most extensively studied latent variable models [Pea94]; see Section 1.4 for the most relevant prior work.

We study both density estimation and parameter estimation for spherical k-GMMs. In density estimation, the goal is to output a hypothesis within small total variation distance from the target. In parameter estimation, the goal is to learn the mean vectors to any desired accuracy—under some necessary assumptions on the pairwise distance between the component means.

Prior to our work, the fastest density estimation algorithm [DK20] had sample and computational complexity $poly(d)(k/\epsilon)^{O(\log^2(k))}$. Here we achieve a polynomial time algorithm under a boundedness assumption on the component means (see Theorem 4.13):

Theorem 1.2 (Density Estimation for Spherical k-GMMs with Bounded Means). There is an algorithm that given $\epsilon > 0$ and $n = \text{poly}(d, k, 1/\epsilon)$ samples from a k-mixture of spherical Gaussians F on \mathbb{R}^d with component means in a ball of radius $O(\sqrt{\log(k)})$, it runs in poly(n, d) time and outputs a hypothesis distribution H such that with high probability $d_{\text{TV}}(H, F) \leq \epsilon$.

This is the first polynomial-time algorithm for density estimation of a natural subclass of spherical mixtures. Theorem 1.2 makes significant progress towards the resolution of a longstanding open problem in the literature, most recently posed in the STOC 2022 talk [Liu22] associated with [LL22].

We now switch our attention to parameter estimation. The state-of-the-art result prior to our work is due to [LL22]. For the case of uniform mixtures¹, they gave a poly(d, k) time algorithm under the assumption that the minimum pairwise mean separation $s = \min_{i \neq j} ||\mu_i - \mu_j||_2$ satisfies $s \gg$ $(\log k)^{1/2+c}$ for any small constant c > 0. The information-theoretically optimal separation under which this task is solvable with poly(d, k) samples is $s \gg (\log k)^{1/2}$ [RV17]. We give a poly $(d, k, 1/\epsilon)$ time parameter estimation algorithm under the optimal separation, with the additional assumption that the largest pairwise distance is comparable to the smallest. Specifically, we establish the following (see Theorem 4.16):

Theorem 1.3 (Clustering Mixtures of Spherical Gaussians under Optimal Separation). Let $F = \sum_{i=1}^{k} (1/k) N(\mu_i, I)$ and $\epsilon > 0$. Suppose that the minimum separation $s := \min_{i \neq j} ||\mu_i - \mu_j||_2$ is at least a sufficiently large constant multiple of $\sqrt{\log(k/\epsilon)}$ and furthermore that $\max_{i\neq j} ||\mu_i - \mu_j||_2 = O(\min_{i\neq j} ||\mu_i - \mu_j||_2)$. There exists an algorithm that given $n = \operatorname{poly}(d, k, 1/\epsilon)$ i.i.d. samples from F, runs in $\operatorname{poly}(n, d)$ time, and outputs estimates $\tilde{\mu}_i$ such that with high probability, for some permutation π of [k], it holds $||\tilde{\mu}_i - \mu_i||_2 \leq \epsilon$, for all $i \in [k]$.

¹For simplicity of exposition, we focus on the uniform case here. Both the results of [LL22] and ours hold for arbitrary weights with modified guarantees.

Prior work [DKS18, HL18, KSS18, DK20] had given quasi-polynomial time parameter estimation algorithms under the optimal separation. In the same context, the polynomial-time algorithm of [LL22] succeeds under the stronger condition that $s \gg \log(k/\epsilon)^{1/2+c}$, for a constant c > 0. We note that they do not require the extra assumption that the largest pairwise distance is comparable to the smallest. To remove this, they developed a complicated recursive clustering argument (see Sections 10 and 11 of [LL22]), which we expect can also be applied mutatis mutandis to our setting.

Learning One-hidden-layer ReLU Networks A one-hidden-layer ReLU network is a function $F : \mathbb{R}^d \to \mathbb{R}$ of the form $F(x) = \sum_{i=1}^k w_i \operatorname{ReLU}(v_i \cdot x)$ for some $w_i \in \mathbb{R}$ and unit vectors $v_i \in \mathbb{R}^d$, where ReLU : $\mathbb{R} \to \mathbb{R}$ is defined as $\operatorname{ReLU}(t) \stackrel{\text{def}}{=} \max\{0, t\}^2$. Here we focus on the natural setting that the weights w_i are positive. Let $\mathcal{C}_{d,k}$ be the corresponding class of functions. The PAC learning problem for $\mathcal{C}_{d,k}$ is the following: The input is a multiset of i.i.d. labeled examples (x, y), where $x \sim N(0, I)$ and y = F(x), for some $F \in \mathcal{C}_{d,k}$. The goal is to output a hypothesis $H : \mathbb{R}^d \to \mathbb{R}$ that with high probability is close to F in L_2 -norm.

A central question in the foundations of deep learning is whether this class (sometimes termed "sums of ReLUs") is learnable in polynomial time, when the marginal distribution is Gaussian. Quoting [CKM21]:

"Ideally, we would like an algorithm with sample complexity and running time that is polynomial in all the relevant parameters. Even for learning arbitrary sums of ReLUs, [...], it remains a major open question to obtain a polynomial-time algorithm [...]."

Prior to our work, the fastest known algorithm [DK20] had super-polynomial complexity, even for constant accuracy. Our main result for this problem is the following (see Theorem 4.7):

Theorem 1.4 (PAC Learning $C_{d,k}$). There is an algorithm that given $\epsilon > 0$, and access to examples (x, F(x)), where $x \sim N(0, I)$ and $F \in C_{d,k}$ is the target function, it runs in time $poly(d, k)2^{poly(1/\epsilon)}$, and outputs a hypothesis $H : \mathbb{R}^d \to \mathbb{R}$ such that with high probability $||H - F||_2 \leq \epsilon ||F||_2$.

Theorem 1.4 makes significant progress towards resolving the aforementioned open problem. Specifically, it gives the first learning algorithm for sums of ReLUs that runs in poly(d, k) time for *sub-constant* values of $\epsilon = o(1)$ —namely up to $\epsilon = 1/\log^c(dk)$, for some constant c > 0. The fastest previous algorithm [DK20] has sample and computational complexity $poly(d)(k/\epsilon)^{O(\log^2(k))}$, i.e., quasi-polynomial in k even for $\epsilon = \Theta(1)$. On the other hand, the running time of the [DK20] algorithm is better than ours for ϵ very small (e.g., 1/poly(d, k)).

Remark 1.5. For one-hidden-layer ReLU networks with general (i.e., positive or negative) weights, the fastest known algorithms take time $(d/\epsilon)^{\text{poly}(k)}$ [CN24, DK24]. Moreover, cryptographic hardness results [LZZ24] rule out the existence of an algorithm with poly(d, k) runtime.

1.2 Main Result: General Implicit Moment Tensor Estimation

All aforementioned learning results are obtained as applications of a general implicit moment tensor computation algorithm that we design, which is the focus of this section. We start by describing our framework, under which such an algorithm is possible, followed by the statement of our main result (Proposition 1.7) and a comparison with the prior work [LL22].

²Our result does not require the assumption that the activation function is a ReLU. The only properties required is that the activation has bounded higher moments and non-vanishing even-degree Fourier coefficients. Our algorithm can be straightforwardly extended to activation functions satisfying these properties.

Sequential Tensor Computations and Implicit Moment Tensor Estimation Roughly speaking, our main result shows that we can efficiently compute a succinct approximate description of tensors of the form $M_m = \sum_{i=1}^k w_i v_i^{\otimes m}$, for non-negative w_i 's, even for large (aka super-constant) values of the order parameter m, assuming there exists a polynomial-size arithmetic circuit whose expected output—under an appropriate samplable distribution—is equal to M_m and whose covariance (under the same distribution) is bounded above.

To formally state our result, we introduce the following key notion:

Definition 1.6 (Sequential Tensor Computation). A sequential tensor computation (STC) is an arithmetic circuit S_t that takes as input \mathcal{I} a set of scalars and vectors in \mathbb{R}^d and outputs an order-t tensor of dimension d, denoted by $S_t(\mathcal{I}) \in (\mathbb{R}^d)^{\otimes t}$, for some $t \in \mathbb{Z}_+$, by applying the following operations:

- (i) Multiplication of a tensor, vector, or scalar by a scalar.
- (ii) Addition of two tensors, vectors or scalars of the same order and dimension.
- (iii) Tensor product of an order-k tensor, for some k < t, by a vector (in the last component) to obtain an order-(k + 1) tensor.

We let the order of S_t be used to denote the order of the output tensor, and the size of S_t to denote the size of the underlying arithmetic circuit.

We note that an STC is a general way to produce higher-order tensors that is compatible with the kinds of computations that we need to perform on them (see Lemmas 3.4 and 3.5) in the context of learning applications. Our main algorithmic result is the following:

Proposition 1.7 (Main Algorithmic Result). Let $w_1, \ldots, w_k \in \mathbb{R}_{\geq 0}$ and $v_1, \ldots, v_k \in \mathbb{R}^d$. Consider the sequence of dimension-d tensors $\{M_t\}$ with M_t order-t, for $t \in \mathbb{Z}_+$, defined by $M_t = \sum_{i=1}^k w_i v_i^{\otimes t}$.

Let $m \in \mathbb{Z}_+$. Suppose that for all positive integers $t \leq 2m$, with t even or equal to m, there exists an STC S_t of order-t and size at most S with the following property: when given as input \mathcal{I} a sample drawn from some efficiently samplable distribution \mathcal{D}_t , the output tensor $S_t(\mathcal{I}), \mathcal{I} \sim \mathcal{D}_t$, has mean M_t and covariance bounded above by V, for some V > 0. Let \mathcal{F}_m be another STC of order-m and size at most S, whose input is partitioned as $(\mathcal{X}, \mathcal{Y})$, with the following properties: when $\mathcal{X} \sim \mathcal{X}$ and $\mathcal{Y} \sim \mathcal{D}'$, where \mathcal{D}' is efficiently samplable, the tensor $\mathcal{F}_m(\mathcal{X}, \mathcal{Y})$ has $\mathbf{E}_{\mathcal{Y} \sim \mathcal{D}'}[\mathcal{F}_m(\mathcal{X}, \mathcal{Y})] = T(\mathcal{X})$ and second moment over both \mathcal{X} and \mathcal{Y} bounded by V.

There is an algorithm that given S, V, k, sample access to $\mathcal{D}_t, \mathcal{D}'$, and $\tau > 0$, it has the following guarantee: on input a single sample $\mathcal{X} \sim X$, the algorithm draws N samples from each of $\mathcal{D}_t, \mathcal{D}'$, runs in poly(N, S, d) time, and with probability $1 - \tau$ (over the samples drawn from $\mathcal{D}_t, \mathcal{D}'$) outputs an approximation \mathcal{A} to $\langle T(\mathcal{X}), M_m \rangle$ such that

$$\mathbf{E}_{\mathcal{X}\sim X}[(\mathcal{A} - \langle T(\mathcal{X}), M_m \rangle)^2] \le \frac{\text{poly}\left(k, m, d, V, 1/\tau, 1 + \sum_{i=1}^k w_i\right) \left(1 + \max_{i \in [k]} \|v_i\|_2\right)^{2m}}{N^{1/2}}$$

In words, we show that, with high probability over the samples drawn from \mathcal{D}_t and \mathcal{D}' , the expected squared error of our estimator over a random choice of $\mathcal{X} \sim X$ is bounded as shown above. In particular, our algorithm can be thought of as taking samples from \mathcal{D}_t and \mathcal{D}' as input and returning an evaluator for a function of \mathcal{X} that with high probability is close in L_2 to the function $\langle T(\mathcal{X}), M_m \rangle$.

Some remarks are in order. First, we did not attempt to optimize the polynomial dependence on the relevant parameters, as this would not affect the qualitative aspects of our learning theory applications. Second, the dependence on the failure probability τ can be made logarithmic, i.e., $\log(1/\tau)$, by running the algorithm $O(\log(1/\tau))$ times and taking a median; such a dependence is not required in our applications. **Comparison to Prior Work** Proposition 1.7 gives an efficient method to approximate higher moment tensors of the form $\sum_{i=1}^{k} w_i v_i^{\otimes t}$ for non-negative w_i 's, where the order t is large, so long as there are "reasonable" unbiased estimators for these moments that can be efficiently sampled. This result broadly generalizes the work by Li and Liu [LL22], who introduced the idea of implicit moment tensor estimation in the context of clustering mixtures of separated spherical Gaussians and, more broadly, separated mixtures of translates of any known Poincare distribution D. Roughly speaking, using the terminology introduced above, they gave an efficient algorithm for implicit moment tensor estimation only for the *specific* STC S_t arising in their clustering application namely, corresponding to the particular expression they used for the adjusted polynomials that come up in their clustering model.

The key conceptual contribution of our work is the realization that the computational task of implicit moment tensor estimation can be solved efficiently in a vastly more general setting, and thus applied to several other learning problems. This answers the main open question posed in [LL22]; see, e.g., [Liu24]. To achieve this goal, we introduce our general framework involving the notion of an STC (Definition 1.6)—which is new to this work. It turns out that this type of *sequence* of tensors arises in a variety of probabilistic models as Fourier or moment tensors. Learning these moment tensors is thus a powerful tool and, as we have shown (Section 1.1), leads to more efficient learning algorithms for these models. Beyond its generality, this broader view leads to quantitative improvements even for the specific clustering application of [LL22], allowing us to cluster mixtures of spherical Gaussians under the optimal separation of $\sqrt{\log(k)}$; see Theorem 1.3 and Remark 4.2.

1.3 Technical Overview

1.3.1 Proof Overview of Proposition 1.7

The naive method to approximate $\langle T(\mathcal{X}), M_m \rangle$ is as follows: Since M_m is the expectation of the output $\mathcal{S}_m(\mathcal{D}_m)$, where \mathcal{D}_m was assumed to be efficiently samplable, we can approximate M_m by drawing sufficiently many samples from \mathcal{D}_m , computing the random output tensors (obtained from \mathcal{S}_m on each sample), and averaging them. Similarly, we can average the output tensors of \mathcal{F}_m , obtained via drawing sufficiently many samples from the relevant input distribution \mathcal{D}' , to get an approximation to $T(\mathcal{X})$. Computing the inner product of these approximations yields an approximation to the desired answer.

This naive approach suffers from two major issues that we need to overcome. The first is statistical and the second is computational.

Statistical Challenges. We start with the statistical aspect. Note that both $T(\mathcal{X})$ and M_m are order-*m* tensors in *d* dimensions, i.e., they lie in a d^m -dimensional space. Thus, in the absence of further structural properties, in order to approximate the means of the sequential tensor computations S_m and \mathcal{F}_m to ℓ_2 -error ϵ , one would need roughly $\Omega(d^m(V/\epsilon^2))$ samples, where *V* is the covariance upper bound. Unfortunately, this sample upper bound is prohibitively large for our learning applications—indeed, in the underlying applications, we need to take $m = \omega(1)$.³

In order to overcome this obstacle, we need a method of reducing the dimension of the underlying space. We can achieve this by considering moment tensors $M_t = \sum_{i=1}^k w_i v_i^{\otimes t}$ of different orders. Suppose, for example, that we knew *exactly* the moment tensor $M_{2m} = \sum_{i=1}^k w_i v_i^{\otimes 2m} =$ $\sum_{i=1}^k w_i v_i^{\otimes m} \otimes v_i^{\otimes m}$. When viewed as a $d^m \times d^m$ matrix, M_{2m} will be rank-k with image (span of

³Concretely, for learning mixtures of k spherical Gaussians one would need $m = \Omega(\log(k))$. For learning sums of ReLU activations, we would need to select $m = \text{poly}(1/\epsilon)$, where ϵ is the desired L_2 -error. Finally, for learning k-mixtures of linear regressions with Gaussian noise of standard deviation σ , one would need $m = \Omega(\log(1/\epsilon)/\sigma^2)$.

its column vectors) W spanned by the $v_i^{\otimes m}$'s. Thus, if we already knew M_{2m} , we could compute its image W and then note that M_m lies in W. Then, instead of computing the full output tensors of the sequential tensor computations \mathcal{S}_m and \mathcal{F}_m —which are objects lying in $(\mathbb{R}^d)^{\otimes m}$ —we could compute their projections onto W. Since dim $(W) \leq k$ (as M_{2m} is rank-k by definition), by standard concentration arguments it follows that $O(kV/\epsilon^2)$ samples suffice to compute these projections to ℓ_2 -error ϵ .

Computational Challenges. The above discussion assumes that we have access to the tensor M_{2m} . This begs the question of how we can compute M_{2m} , even approximately. It turns out that this can be achieved via some form of bootstrapping. In particular, suppose that we have computed (an approximation to) M_{2t} , for some $t \leq m$. We can write $M_{2t} = \sum_{i=1}^{k} w_i v_i^{\otimes t} \otimes v_i^{\otimes t}$. When viewed as a $d^t \times d^t$ matrix, M_{2t} has image spanned by the $v_i^{\otimes t}$'s. Since the image is at most k-dimensional, if W_t is the span of the top-k singular vectors of (our approximation to) M_{2t} , then the $v_i^{\otimes t}$'s will all (approximately) lie in W_t . This means that M_{2t+2} will (approximately) lie in the space

$$W_t \otimes \mathbb{R}^d \otimes W_t \otimes \mathbb{R}^d$$

As this space has dimension $\dim(W_t)^2 d^2 = O(k^2 d^2)$, this allows a relatively fast approximation to M_{2t+2} . Thus, an approximation to M_{2t} allows us to compute an approximation to M_{2t+2} . A careful analysis of the error rates involved resolves our statistical issues.

While the aforementioned ideas can be used to overcome sample complexity considerations, significant computational challenges remain. In particular, all of the vectors/tensors in question still lie in d^t dimensional spaces; so, even writing them down explicitly will take $\Omega(d^t)$ time. We circumvent this obstacle by devising a succinct representation of the relevant vector spaces that suffices for the purpose of the underlying computations. In particular, by the above discussion, we can express W_{t+1} (the span of the top-k singular vectors of our approximation to $M_{2(t+1)}$) as a subspace of $W_t \otimes \mathbb{R}^d$. By associating this (k-dimensional) subspace with \mathbb{R}^k (by picking an appropriate orthonormal basis/isometry), we can then associate W_{t+2} with a subspace of $\mathbb{R}^k \otimes$ \mathbb{R}^d ; and so on. In particular, for each $t \leq m$, we will have a map $\Phi_t : (\mathbb{R}^d)^{\otimes t} \to \mathbb{R}^k$ given by the projection onto W_t followed by an isometry. These maps can then be efficiently constructed recursively from each other.

1.3.2 From Proposition 1.7 to Learning Latent-variable Models

In this subsection, we sketch how our main result is leveraged to derive our learning applications. We start with the problem of learning sums of ReLUs, followed by our density estimation algorithms for MLRs and GMMs, and concluding with our application on clustering GMMs.

PAC Learning Sums of ReLUs Recall that the objective is to approximate a function of the form $F(x) = \sum_{i=1}^{k} w_i \text{ReLU}(v_i \cdot x)$, where $w_i \ge 0$, given access to random labeled examples (x, F(x)) with $x \sim N(0, I)$. We proceed by attempting to learn the low degree part of its Fourier transform. In particular, we have that

$$F(x) = \sum_{n=0}^{\infty} \langle T_n, H_n(x) \rangle$$

where H_n is the normalized Hermite tensor (Definition 2.2) and $T_n = \mathbf{E}[F(G)H_n(G)]$. Our approach relies on approximating F by the sum of the first $\text{poly}(1/\epsilon)$ many terms in the sum above. To do this, we note that the tensors T_n are moment tensors of approximately the kind desired. Namely, it can be shown (see Lemma 4.8) that

$$T_n = c_n \sum_{i=1}^k w_i v_i^{\otimes n} ,$$

for some known constants c_n (which are critically bounded away from 0 when n is even). Given this, we can use Proposition 1.7 directly to approximate $\langle T_n, H_n(x) \rangle$ in time roughly poly $(k, d)2^{O(n)}$. Doing this for all n up to degree poly $(1/\epsilon)$ yields the desired approximation (proof of Theorem 4.7).

Density Estimation for MLRs and spherical GMMs Our algorithms for density estimation, both for mixtures of spherical Gaussians and mixtures of linear regressions, proceed similarly to the function approximation one above. In particular, in order to do density estimation for a distribution X, our goal will be to approximate the function F(x) := X(x)/G(x), where G is the standard Gaussian. Once we have this, we can simulate X via rejection sampling, by producing a sample $x \sim G$ and accepting with probability proportional to F(x).

As above, F(X) has a Taylor expansion given by

$$F(x) = \sum_{n=0}^{\infty} \langle T_n, H_n(x) \rangle ,$$

where $T_n = \mathbf{E}[H_n(X)]$. This in turn needs to be related to an appropriate moment tensor. For Gaussian mixtures, T_n is already $\sum_{i=1}^k w_i \mu_i^{\otimes n}$. For mixtures of linear regressions on the other hand, it needs to be computed indirectly as $\sum_{i=1}^k w_i \beta_i^{\otimes n} = \mathbf{E}[(y^n/\sqrt{n!})H_n(X)]$, and the Fourier coefficient is given by (see Lemma 4.20)

$$T_n = \frac{(n-1)!!}{\sqrt{n!}} \sum_{i=1}^k w_i \operatorname{Sym}\left(\begin{bmatrix} 0 & \beta_i \\ \beta_i^T & \|\beta_i\|_2^2 + \sigma^2 - 1 \end{bmatrix}^{\otimes n/2} \right)$$

The inner product of this quantity with $H_n(x)$ can be obtained indirectly by taking appropriate inner products with the moment tensor.

Clustering Spherical GMMs The basic idea for our clustering application mirrors the approach of [LL22]: given two samples x and x' from our mixture $X = \sum_{i=1}^{k} w_i N(\mu_i, I)$, we want to reliably determine whether or not x and x' are from the same component or different ones. If we can cluster samples by component, we can then easily learn each component. We do this by considering the distribution $X' = (Y - Y')/\sqrt{2}$, where Y and Y' are independent copies of X. Here X' is another mixture of Gaussians with means $\mu_i - \mu_j$. We can determine whether x and x' are in the same component by considering the size of the inner product of $(x - x')^{\otimes t}$ (actually in our case $H_t((x - x')/\sqrt{2})$) with the moment tensor of X', namely $\sum_{i,j} w_i w_j (\mu_i - \mu_j)^{\otimes t} = \mathbf{E}[H_t(X')]$; indeed, this inner product will be large with high probability if and only if x and x' come from different components of the mixture.

While the above is essentially the strategy employed by [LL22], we devise a more efficient sequential tensor computation for estimating the values of Hermite tensors (see Definition 4.1 and Lemma 4.3). In particular, this allows us to efficiently consider order $t = \Theta(\log(k))$ tensors, while [LL22] could only consider order up to $O(\log(k)/\log\log(k))$. This quantitative improvement is responsible for our improvement in the separation parameter.

1.4 Prior and Related Work

Here we record the most relevant prior work on the learning problems we study.

Learning Mixtures of Linear Regressions A line of prior work on learning MLRs focused on the parameter estimation problem. Specifically, a sequence of paper (see, e.g., [ZJD16, LL18, KC19] and references therein) analyzed non-convex methods, including expectation maximization and alternating minimization. These works established *local* convergence guarantees: Given a sufficiently accurate solution (warm start), these non-convex methods can efficiently boost this to a solution with arbitrarily high accuracy. The focus of our algorithmic results in this section is to provide such a warm start, which captures the complexity of the problem. We note that the local convergence result of [LL18] applies for the noiseless case, while the result of [KC19] can handle non-trivial regression noise when the weights of the unknown mixture are known.

The prior works most closely related to ours are [LL18, CLS20, DK20]. The work of [LL18] studied the noiseless setting (corresponding to $\sigma = 0$) and provided an algorithm with sample complexity and running time scaling exponentially with k. Subsequently, [CLS20] gave an algorithm whose sample complexity and running time scales with $\exp(\tilde{O}(k^{1/2}))$, i.e., sub-exponential in k. In more detail, for the noisy case, when $\sigma = O(\epsilon)$ and the weights are uniform, the algorithm of [CLS20] has sample complexity and runtime of the form $\operatorname{poly}(dk/(\epsilon\Delta))(k/\epsilon)^{\tilde{O}(k^{1/2}/\Delta^2)}$, where Δ is the minimum pairwise separation between the regressors. The fastest previously known algorithm [DK20] for both parameter and density estimation has sample and computational complexity scaling quasi-polynomially in k. Specifically, for density estimation the latter algorithm has complexity poly $(d, (k/\epsilon)^{\log^2(k)})$.

Learning Mixtures of Spherical Gaussians Here we survey the most relevant prior work on learning this distribution family both for density estimation and parameter estimation.

Density estimation for mixtures of high-dimensional spherical Gaussians has been studied in a series of works [FOS06, MV10, SOAJ14, HP15, DKK⁺16, LS17, ABH⁺18]. The sample complexity of the problem for k-mixtures on \mathbb{R}^d , for variation distance error ϵ , is $\tilde{\Theta}(dk/\epsilon^2)$ [ABH⁺18]. Unfortunately, until fairly recently, all known algorithms had running times exponential in number k; see [SOAJ14] and references therein. The fastest density estimation algorithm prior to our work is from [DK20] and has complexity poly $(d)(k/\epsilon)^{O(\log^2(k))}$.

In the related task of parameter estimation, the goal is to approximate the parameters of the components (i.e., mixture weights and component means). This task requires further assumptions, so that it is solvable with polynomial sample complexity (even information-theoretically). The typical assumption in the literature is that pairwise separation between the component means is bounded below by some parameter $\Delta > 0$. A long line of work since the late 90s steadily improved the separation requirement [Das99, AK01, VW02, AM05, KSV08, BV08, RV17, HL18, KSS18, DKS18, LL22]. The state-of-the-art prior to this paper is the work of [LL22] that gave a poly $(d, k, 1/\epsilon)$ time algorithm under near-optimal separation of $\Theta((\log(k/\epsilon))^{1/2+c})$. Finally, we note that the assumption on spherical component covariances is crucial for the latter result (and our Theorem 1.3). Namely, a departure from sphericity leads to information-computation tradeoffs: [DKPZ23] gave evidence that quasi-polynomial runtime is inherent for any polylog(k) separation, even if the Gaussian components have the same bounded covariance.

Learning One-hidden-layer ReLU Networks There has been a recent explosion of research on provable algorithms for learning neural networks in various settings, see, e.g., [JSA15, SJA16, DFS16, ZLJ16, ZSJ⁺17, GLM18, GKLW19, BJW19, GKKT17, GK19, VW19, DKKZ20, GGJ⁺20, DK20, CKM21, CGKM22, CDG⁺23, DK24] for some works on the topic. Many of these works focused on parameter learning—the problem of recovering the weight matrix of the data generating neural network. PAC learning of simple classes of networks has been studied as well [GKKT17, GK19, VW19, DKKZ20, CGKM22, CDG⁺23, DK24].

The work of [GLM18] studies the parameter learning of positive linear combinations of ReLUs under the Gaussian distribution. It is shown in [GLM18] that the parameters can be efficiently approximated, if the weight matrix is full-rank with bounded condition number. The complexity of their algorithm scales polynomially with the condition number. [BJW19, GKLW19] obtained efficient parameter learners for vector-valued depth-2 ReLU networks under the Gaussian distribution. Similarly, the algorithms in these works have sample complexity and running time scaling polynomially with the condition number.

In contrast to parameter estimation, PAC learning these function classes does not require any assumptions on the structure of the weight matrix. The PAC learning problem for this class is information-theoretically solvable with polynomially many samples. The question is whether a computationally efficient algorithm exists. Until recently, the problem of PAC learning positive linear combinations of ReLUs For the task of positive linear combinations of ReLUs studied in this work, [DKKZ20] gave a learner with sample complexity $poly(dk/\epsilon)$ and computational complexity $poly(dk/\epsilon) + (k/\epsilon)^{O(k^2)}$. This runtime bound was improved to $poly(d)(k/\epsilon)^{O(\log^2(k))}$ in [DK20], which was the state-of-the-art prior to our work.

1.5 Organization

The structure of the paper is as follows: In Section 2, we provide the necessary definitions and technical facts. In Section 3, we prove our main result. Section 4 presents our learning algorithms for sums of ReLUs, mixtures of spherical Gaussians, and mixtures of linear regressions.

2 Preliminaries

Notation For $n \in \mathbb{Z}_+$, we denote by [n] the set $\{1, 2, \ldots, n\}$. For a vector $v \in \mathbb{R}^d$, let $||v||_2$ denote its Euclidean norm. We denote by $x \cdot y$ the standard inner product between $x, y \in \mathbb{R}^d$. We will denote by δ_0 the Dirac delta function and by $\delta_{i,j}$ the Kronecker delta. Throughout the paper, we let \otimes denote the tensor/Kronecker product. For a vector $x \in \mathbb{R}^d$, we denote by $x^{\otimes m}$ the *m*-th order tensor power of x.

Let V be an inner product space. If A and B are elements of $V^{\otimes t}$ for some $t \in \mathbb{Z}_+$, then we use $\langle A, B \rangle$ to denote the inner product of A and B induced by the inner product on V. We also use $||A||_2 = \langle A, A \rangle^{1/2}$ for the corresponding ℓ_2 -norm.

A tensor A of order t and dimensions $(d_i)_{i \in [t]}$ is a multilinear map defined by a t-dimensional array with real entries A_{α} , where $\alpha = (\alpha_1, \ldots, \alpha_t)$ with $\alpha_i \in [d_i]$. The covariance of an order-t, dimension d tensor-valued random variable $T \in (\mathbb{R}^d)^{\otimes t}$ is defined as the tensor in $(\mathbb{R}^d)^{\otimes 2t}$ defined as $\mathbf{Cov}[T] := \mathbf{E}[T \otimes T] - \mathbf{E}[T] \otimes \mathbf{E}[T]$, where the expectation operation is applied entry-wise. Similarly, the corresponding "second moment" is the tensor $\mathbf{E}[T \otimes T]$. For an order-t, dimension-d tensor W, consider the real-valued random variable $\langle W, T \rangle$ and its corresponding variance, $\mathbf{Var}[\langle W, T \rangle]$. We say that $\mathbf{Cov}[T]$ is bounded above by C, for some C > 0, if $\max_{\|W\|_2=1} \mathbf{Var}[\langle W, T \rangle] \leq C$, i.e., if the variance of T in any (normalized) direction is at most C. The analogous definition applies for the second moment $\mathbf{E}[T \otimes T]$.

For a space W, we use I_W for the identity matrix on W. If $W = \mathbb{R}^d$, we use the notation I_d .

We will denote by $N(0, I_d)$ the *d*-dimensional Gaussian distribution with zero mean and identity covariance; we will use N(0, I) when the underlying dimension will be clear from the context.

We will use N(0,1) for the univariate case. For a random variable X and $p \ge 1$, we will use $||X||_p \stackrel{\text{def}}{=} \mathbf{E}[|X|^p]^{1/p}$ to denote its L_p -norm.

Hermite Analysis and Concentration Consider $L_2(\mathbb{R}^d, N(0, I))$, the vector space of all functions $f : \mathbb{R}^d \to \mathbb{R}$ such that $\mathbf{E}_{x \sim N(0,I)}[f(x)^2] < \infty$. This is an inner product space under the inner product $\langle f, g \rangle = \mathbf{E}_{x \sim N(0,I)}[f(x)g(x)]$. This inner product space has a complete orthogonal basis given by the *Hermite polynomials*. In the univariate case, we will work with normalized Hermite polynomials defined below.

Definition 2.1 (Normalized Probabilist's Hermite Polynomial). For $k \in \mathbb{N}$, the k-th probabilist's Hermite polynomial $He_k : \mathbb{R} \to \mathbb{R}$ is defined as $He_k(t) = (-1)^k e^{t^2/2} \cdot \frac{d^k}{dt^k} e^{-t^2/2}$. We define the k-th normalized probabilist's Hermite polynomial $h_k : \mathbb{R} \to \mathbb{R}$ as $h_k(t) = He_k(t)/\sqrt{k!}$.

Note that for $G \sim N(0,1)$ we have $\mathbf{E}[h_n(G)h_m(G)] = \delta_{n,m}$. We will use multivariate Hermite polynomials in the form of Hermite tensors. We define the normalized Hermite tensor as follows, in terms of Einstein summation notation.

Definition 2.2 (Normalized Hermite Tensor). For $k \in \mathbb{N}$ and $x \in V$ for some inner produce space V, we define the k-th Hermite tensor as

$$(H_k^{(V)}(x))_{i_1,i_2,\dots,i_k} := \frac{1}{\sqrt{k!}} \sum_{\substack{\text{Partitions } P \text{ of } [k]\\ \text{into sets of size 1 and 2}}} \bigotimes_{\{a,b\} \in P} (-I_{i_a,i_b}) \bigotimes_{\{c\} \in P} x_{i_c} ,$$

where I above denotes the identity matrix over V. Furthermore, if $V = \mathbb{R}^d$, we will often omit the superscript and simply write $H_k(x)$.

We will require a few properties that follow from this definition. First, note that if V is a subspace of W, then $H_k^{(V)}(\operatorname{Proj}_V(x)) = \operatorname{Proj}_V^{\otimes k} H_k^{(W)}(x)$. Applying this when V is the one-dimensional subspace spanned by a unit vector v gives that $\langle H_k(x), v^{\otimes k} \rangle = h_k(v \cdot x)$. We will also need to know that the entries of $H_k(x)$ form a useful Fourier basis of $L^2(\mathbb{R}^d, N(0, I))$. In particular, for non-negative integers m and k, we have that $\mathbf{E}_{x \sim N(0,I)}[H_k(x) \otimes H_m(x)]$ is 0 if $m \neq k$ and $\operatorname{Sym}_k(I_{d^k})$, if m = k, where Sym_k is the symmetrization operation over the first k coordinates. From this we conclude that if T is a symmetric k-tensor, then $\mathbf{E}_{x \sim N(0,I)}[\langle H_k(x), T \rangle H_m(x)]$ is 0 if $m \neq k$ and T if m = k.

We will also require the following fact, where Sym denotes the symmetrization operator that averages a tensor over all permutations of its entries.

Fact 2.3. We have that $\mathbf{E}_{X \sim N(\mu,I)}[H_n(X)] = \mu^{\otimes n} / \sqrt{n!}$ and $\mathbf{Cov}_{X \sim N(0,I)}[H_n(X)] = \mathrm{Sym}(I_{d^n}).$

For a polynomial $p : \mathbb{R}^d \to \mathbb{R}$, we will use $\|p\|_r \stackrel{\text{def}}{=} \mathbf{E}_{x \sim N(0,I)}[|p(x)|^r]^{1/r}$, for $r \geq 1$. We recall the following well-known hypercontractive inequality [Bon70, Gro75]:

Fact 2.4. Let $p: \mathbb{R}^d \to \mathbb{R}$ be a degree-k polynomial and q > 2. Then $\|p\|_q \leq (q-1)^{k/2} \|p\|_2$.

3 Main Result: Proof of Proposition 1.7

Before we describe and analyze our algorithm, we introduce some necessary technical tools.

3.1 Recursive Pseudo-projections and their Properties

We start by defining the notion of a pseudo-projection.

Definition 3.1 (Pseudo-projection). A linear transformation $A : X \to W$ between finite dimensional inner product spaces X, W is called a *pseudo-projection* if $AA^{\top} = I_W$.

The following lemma establishes two basic properties of pseudo-projections that we will require.

Lemma 3.2 (Properties of Pseudo-projections). The following properties hold:

- 1. The composition of pseudo-projections is a pseudo-projection.
- 2. A pseudo-projection $A: X \to W$ can be written as a projection $P: X \to U$, for some subspace U of X, composed with an isometry from U to W.

Proof. Property (1) follows immediately from the definition. For property (2), note that the right singular values of A are equal to 1, and its left singular values are 1 and 0. Therefore, $A^{\top}A$ is a projection onto some subspace U of X of the same dimension as W. Since A^{\top} maps W to U and preserves norms, it is an isometry. Since A, mapping U to W, is its inverse, it follows that $A: U \to W$ (the restriction of A on U) is also an isometry. Note that $A: X \to W$ can be written as $A = A(A^{\top}A) = (AA^{\top})A$, which is the composition of a projection onto U and an isometry. \Box

In order to describe the *recursive* projections, we introduce the following definition:

Definition 3.3 (Recursive Pseudo-projection). For $d, t, n \in \mathbb{Z}_+$, a (d, t, n)-recursive pseudo-projection is a pseudo-projection $\Phi : (\mathbb{R}^d)^{\otimes t} \to \mathbb{R}^n$, defined as follows. There exist positive integers n_0, n_1, \ldots, n_t with $n_0 = 1$ and $n_t = n$ and base pseudo-projections $\phi_i : \mathbb{R}^{n_{i-1}} \otimes \mathbb{R}^d \to \mathbb{R}^{n_i}, i \in [t]$, such that the following holds: We define $\Phi : \mathbb{R}^1 \otimes (\mathbb{R}^d)^{\otimes t} \to \mathbb{R}^n$ by starting with a tensor T_0 in $\mathbb{R}^{n_0} \otimes (\mathbb{R}^d)^{\otimes t}$, applying ϕ_1 to the first two components of T_0 to get a tensor T_1 in $\mathbb{R}^{n_1} \otimes (\mathbb{R}^d)^{\otimes (t-1)}$, then applying ϕ_2 to the first two components of T_1 to get a tensor T_2 in $\mathbb{R}^{n_2} \otimes (\mathbb{R}^d)^{\otimes (t-2)}$; and more generally applying $\phi_i, i \in [t]$, to the first two components of T_{i-1} , which is a tensor in $\mathbb{R}^{n_{i-1}} \otimes (\mathbb{R}^d)^{\otimes (t-i+1)}$, to get a tensor T_i in $\mathbb{R}^{n_i} \otimes (\mathbb{R}^d)^{\otimes (t-i)}$. We say that the recursive pseudo-projection Φ has order tand size max_i n_i .

It is easy to verify that a recursive pseudo-projection is itself a pseudo-projection. We will require a few basic properties of recursive pseudo-projections. First, we show that a recursive pseudo-projection can efficiently be applied to the output of a sequential tensor computation:

Lemma 3.4 (Efficient Computation). Let $\Phi : (\mathbb{R}^d)^{\otimes t} \to \mathbb{R}^n$ be an order-t recursive pseudoprojection whose defining pseudo-projections ϕ_i are given explicitly as matrices. Let S be a sequential tensor computation returning an order-t and dimension-d tensor. Then there is an algorithm that given input vectors v_i for S returns $\Phi(S(\{v_i\}))$ and runs in time poly $(d, \text{size}(\Phi), \text{size}(S))$.

Proof. At any point in the circuit defining S if there is an order-*s* tensor, we can compute the first *s* steps of Φ applied to that value. The goal is to for every wire in the circuit computing S to compute the value of the appropriate iterate of the operation defining Φ applied to the value being carried by that wire in the circuit. We claim that this can be computed efficiently via dynamic programming. We start with the wires at the beginning of the circuit and work our way towards the end. It suffices to verify that if we know the appropriate values for the input wires, we can efficiently compute the corresponding value of the output wire. This involves the operations of tensor summation, multiplication of a tensor by a scalar sum, and tensor product of tensor with a vector in the last coordinate. It is straightforward to see that any of the allowed operations in S are compatible with this computation and that it runs in polynomial time in the relevant parameters.

We also need to be able to efficiently take tensor products of recursive pseudo-projections:

Lemma 3.5 (Efficient Tensorization). Given recursive pseudo-projections $\Phi : (\mathbb{R}^d)^{\otimes t} \to \mathbb{R}^n$ and $\Theta : (\mathbb{R}^d)^{\otimes s} \to \mathbb{R}^m$ along with their defining pseudo-projections ϕ_i, θ_j , given explicitly as matrices, one can efficiently construct the recursive pseudo-projection $\Phi \otimes \Theta : (\mathbb{R}^d)^{\otimes (t+s)} \to \mathbb{R}^{nm}$ of size $\max(\operatorname{size}(\Phi), n \operatorname{size}(\Theta))$.

Proof. The proof follows by using as the base pseudo-projections $\phi_1, \phi_2, \ldots, \phi_t, I_n \otimes \theta_1, \ldots, I_n \otimes \theta_s$.

3.2 Algorithm and Analysis

We are now ready to present the pseudo-code of the algorithm establishing Proposition 1.7.

```
Algorithm Implicit-Moment-Estimation
```

- 1. Let N be a positive integer quantifying the number of samples drawn from the relevant distributions in each step.
- 2. Let A_1 be the average of N runs of S_2 on N independent samples drawn from \mathcal{D}_2 , thought of as a $d \times d$ matrix.
- 3. Let W_1 be the span of the top-k singular vectors of A_1 . Let $\phi_1 : \mathbb{R}^1 \otimes \mathbb{R}^d \to \mathbb{R}^k$ be the projection of \mathbb{R}^d onto W_1 composed with an isometry from W_1 to \mathbb{R}^k .
- 4. For r = 1 to m 1:
 - (a) Let $\Phi_r : (\mathbb{R}^d)^{\otimes r} \to \mathbb{R}^k$ be the order-*r* recursive pseudo-projection given by ϕ_1, \ldots, ϕ_r .
 - (b) Let $\Phi'_r : (\mathbb{R}^d)^{\otimes (2r+2)} \to \mathbb{R}^{k^2 d^2}$ be the recursive pseudo-projection $\Phi_r \otimes I_d \otimes \Phi_r \otimes I_d$.
 - (c) Let A_{r+1} be the average of N copies of Φ'_r applied to the outputs of S_{2r+2} on N independent samples drawn from \mathcal{D}_{2r+2} , using Lemma 3.4 to compute efficiently, thought of as a $dk \times dk$ matrix.
 - (d) Let $W_{r+1} \subset \mathbb{R}^{dk}$ be the span of the top-k singular vectors of A_{r+1} .
 - (e) Let ϕ_{r+1} be the composition of the projection of $\mathbb{R}^k \otimes \mathbb{R}^d \to W_{r+1}$ with an isometry from W_{r+1} to \mathbb{R}^k .
- 5. Let $\Phi_m : (\mathbb{R}^d)^{\otimes m} \to \mathbb{R}^k$ be the recursive pseudo-projection given by ϕ_1, \ldots, ϕ_m .
- 6. Let A be the average of N copies of Φ_m applied to the output of S_m on N independent samples drawn from \mathcal{D}_m , using Lemma 3.4 to compute efficiently.
- 7. Let $B(\mathcal{X})$ be the average of N copies of Φ_m applied to the output of \mathcal{F}_m on the single sample \mathcal{X} and N independent samples from \mathcal{D}' , using Lemma 3.4 to compute efficiently.

Return $\langle A, B(\mathcal{X}) \rangle$.

Proof of Proposition 1.7 By the definition of the A_r 's in the algorithm pseudocode we have that (i) A_1 will be close to M_2 , and (ii) A_{r+1} , $r \in [m-1]$, will be close to $\Phi'_r(M_{2r+2})$ with high probability. Specifically, with probability at least $1 - \tau/4$ over the samples drawn from the \mathcal{D}_t 's in

the various steps, we will have $||A_1 - M_2||_2 \le \delta$ and $||A_{r+1} - \Phi'_r(M_{2r+2})||_2 \le \delta$, for all $r \in [m-1]$, where

$$\delta \stackrel{\text{def}}{=} O\left(\frac{m \, d \, k \, V}{\sqrt{N \, \tau}}\right) \,. \tag{1}$$

We will condition on the event that all these approximations hold in the subsequent analysis.

Let $\Phi_r : (\mathbb{R}^d)^{\otimes r} \to \mathbb{R}^k$ be the recursive pseudo-projection given by ϕ_1, \ldots, ϕ_r . By the second statement of Lemma 3.2, Φ_r is the composition of the projection of $(\mathbb{R}^d)^{\otimes r}$ onto some subspace U_r of $(\mathbb{R}^d)^{\otimes r}$ composed with an isometry.

For $i \in [k]$ and $r \in [m]$, we will denote $x_{i,r} \stackrel{\text{def}}{=} \sqrt{w_i} v_i^{\otimes r}$. With this notation, we can write $M_{2r} = \sum_{i=1}^k x_{i,r}^{\otimes 2}$. Let

$$\eta_r \stackrel{\text{def}}{=} \max_{i \in [k]} \operatorname{dist}(x_{i,r}, U_r) = \max_{i \in [k]} \left\| x_{i,r} - \operatorname{proj}_{U_r}(x_{i,r}) \right\|_2,$$

where we recall that U_r is the subspace of $(\mathbb{R}^d)^{\otimes r}$ corresponding to Φ_r .

We will establish the following lemma:

Lemma 3.6. For all $r \in [m-1]$, it holds $\eta_{r+1} \leq \eta_r \max_i ||v_i||_2 + O(\sqrt{\delta})$.

Proof. By the definition of Φ'_r and M_{2r+2} , we can write:

$$\Phi'_r(M_{2r+2}) = \sum_{i=1}^k (\underbrace{\Phi_r(x_{i,r}) \otimes v_i}_{y_{i,r}})^{\otimes 2} = \sum_{i=1}^k y_{i,r}^{\otimes 2}.$$

Moreover, we have that $(\Phi_r \otimes I_d)^{\top} y_{i,r} = \Phi_r^{\top} (\Phi_r(x_{i,r})) \otimes v_i = \operatorname{Proj}_{U_r}(x_{i,r}) \otimes v_i$. Therefore,

$$\left\| (\Phi_r \otimes I_d)^\top y_{i,r} - x_{i,r+1} \right\|_2 \le \eta_r \|v_i\|_2 .$$
 (2)

There are three relevant objects for the analysis: M_{2r+2} ,

$$M^* \stackrel{\text{def}}{=} (\Phi'_r)^\top \Phi'_r M_{2r+2} = \sum_{i=1}^k (\Phi'_r)^\top \Phi'_r (y_{i,r}^{\otimes 2}) ,$$

and $A^* \stackrel{\text{def}}{=} (\Phi'_r)^\top A_{r+1}$.

We note that

$$||A^* - M^*||_2 = ||A_{r+1} - \Phi'_r M_{2r+2}||_2 \le \delta$$

Recall that $\eta_{r+1} = \max_{i \in [k]} \operatorname{dist}(x_{i,r+1}, U_{r+1})$. For all $i \in [k]$,

$$\operatorname{dist}(x_{i,r+1}, U_{r+1}) \le \|x_{i,r+1} - (\Phi_r \otimes I_d)^\top y_{i,r}\|_2 + \operatorname{dist}((\Phi_r \otimes I_d)^\top y_{i,r}, U_{r+1}) .$$

It follows from (2) that

$$|x_{i,r+1} - (\Phi_r \otimes I_d)^{\top} y_{i,r}||_2 \le \eta_r ||v_i||_2 .$$

Moreover, note that dist $((\Phi_r \otimes I_d)^\top y_{i,r}, U_{r+1})$ is the distance from $(\Phi_r \otimes I_d)^\top y_{i,r}$ to the span of the top-k singular vectors of A^* . Since A^* is $O(\delta)$ -close to $M^* = \sum_{i=1}^k (\Phi'_r)^\top y_{i,r}^{\otimes 2}$ (which is rank-k), we can show that dist $((\Phi_r \otimes I_d)^\top y_{i,r}, U_{r+1}) = O(\sqrt{\delta})$. We prove this as follows. Let $y := (\Phi_r \otimes I_d)^\top y_{i,r} = u + v$ with $u \in U_{r+1}$ and v orthogonal. The distance in question is now just $\|v\|_2$. Note that $v^\top M^* v \ge \langle v, y \rangle^2 \ge \|v\|_2^4$. On the other hand $v^\top A^* v$ must be relatively small. This is because v is orthogonal to the top-k singular vectors of A^* , and so it is at most $||v||_2^2 \lambda_{k+1}$, where λ_{k+1} is the (k+1)-st singular value. We have that

$$\lambda_{k+1} \le \sup_{w \perp (\Phi_r \otimes I_d)^\top y_{i,r}, 1 \le i \le k, \|w\|_2 = 1} w^\top A^* w = \sup_{w \perp (\Phi_r \otimes I_d)^\top y_{i,r}, 1 \le i \le k, \|w\|_2 = 1} w^\top M^* w + O(\delta) = O(\delta) .$$

So $v^{\top} M^* v \ge \|v\|_2^4$, and $v^{\top} A^* v = O(\delta) \|v\|_2^2$, which gives that the difference is

$$v^{\top} (A^* - M^*) v = O(\delta) \|v\|_2^2$$

Putting these together implies that $||v||_2 = O(\sqrt{\delta})$.

Thus, we have that $\eta_{r+1} = O(\eta_r \max_i ||v_i||_2 + \sqrt{\delta})$, proving Lemma 3.6.

Lemma 3.6 allows us to establish an appropriate upper bound on η_m , which suffices to show that $\Phi_m^{\top} A$ is close to M_m with high probability. This in turns gives that $\langle A, B(X) \rangle$ close to $\langle M_m, T(X) \rangle$ in the mean squared sense. We present the detailed argument below.

We start by using Lemma 3.6 to show an upper bound of η_m . To do this, we need an upper bound on $\eta_1 = \max_{i \in [k]} \operatorname{dist}(x_{i,1}, U_1)$. Recall that $x_{i,1} = \sqrt{w_i}v_i$ and note that by construction it holds $U_1 = W_1$, where W_1 is the span of the top-k singular vectors of A_1 .

We will use the fact that $||A_1 - M_2||_2 \leq \delta$. In the limit, when $\delta \to 0$, we have that $A_1 = M_2 = \sum_{i=1}^k x_{i,1} x_{i,1}^\top$, in which case W_1 is the span of the $x_{i,1}$'s and $\eta_1 = 0$. For the case of $\delta > 0$, a standard argument using Weyl's inequality (similar to the one used in the proof of Lemma 3.6) shows that $||x_{i,1} - \operatorname{proj}_{W_1}(x_{i,1})||_2 = O(\sqrt{\delta})$, which gives that $\eta_1 = O(\sqrt{\delta})$.

Unrolling the recursion, and using the fact that $\eta_1 = O(\sqrt{\delta})$, we obtain

$$\eta_m = O(m\sqrt{\delta}) \max\left\{1, \left(\max_{i \in [k]} \|v_i\|_2\right)^{m-1}\right\}$$
(3)

We now write Φ_m as the composition of a projection \mathbb{P}_m onto U_m composed with an isometry $\mathbb{R}_m : U_m \to \mathbb{R}^k$. We establish the following claim:

Claim 3.7. We have that $||M_m - P_m(M_m)||_2 \le s \stackrel{\text{def}}{=} \eta_m \sqrt{k} \left(\sum_{i=1}^k w_i\right)^{1/2}$.

Proof. Since $M_m = \sum_{i=1}^k \sqrt{w_i} x_{i,m}$, by linearity of P_m it follows that $P_m(M_m) = \sum_{i=1}^k \sqrt{w_i} P_m(x_{i,m})$. Therefore,

$$\|M_m - \mathcal{P}_m(M_m)\|_2 = \sum_{i=1}^k \sqrt{w_i} \|x_{i,m} - \mathcal{P}_m(x_{i,m})\|_2$$
$$\leq \eta_m \sum_{i=1}^k \sqrt{w_i}$$
$$\leq \eta_m \sqrt{k} \left(\sum_{i=1}^k w_i\right)^{1/2},$$

where the first inequality follows from the definition of η_m and the second is Cauchy-Schwarz.

By construction, A is the average of N copies of $R_m(P_m(\mathcal{S}_m))$ on independent samples from the distribution \mathcal{D}_m . By assumption, $P_m(\mathcal{S}_m)$ has mean $P_m(M_m)$ and covariance bounded by V

(since P_m is a projection). Given that these random variables lie in U_m , which is k-dimensional, we have that with probability $1 - \tau/4$ over samples from \mathcal{D}_m it holds

$$\|Average(\mathbf{P}_m(\mathcal{S}_m)) - \mathbf{P}_m(M_m)\|_2 = O\left(\sqrt{\frac{kV}{N\tau}}\right)$$

Since R_m is an isometry, we also get that with probability $1 - \tau/4$ over samples from \mathcal{D}_m

$$||A - \mathbf{R}_m(\mathbf{P}_m(M_m))||_2 = O\left(\sqrt{\frac{kV}{N\tau}}\right)$$

By similar logic, with probability $1 - \tau/4$ over samples from \mathcal{D}' , we have that

$$\mathbf{E}_{\mathcal{X}\sim X}[\|B(\mathcal{X}) - \mathbf{R}_m(\mathbf{P}_m(T(\mathcal{X})))\|_2^2]^{1/2} = O\left(\sqrt{\frac{kV}{N\tau}}\right) \ .$$

Using the above, with probability at least $1 - \tau$ over the samples drawn from \mathcal{D} and \mathcal{D}' , we obtain the following chain of (in)equalities.

$$\langle A, B(\mathcal{X}) \rangle$$

$$= \langle \mathcal{R}_m(\mathcal{P}_m(T(\mathcal{X}))), \mathcal{R}_m(\mathcal{P}_m(M_m)) \rangle + O\left(\sqrt{kV/(N\tau)}\right) (\|\mathcal{P}_m(T(\mathcal{X}))\|_2 + \|\mathcal{P}_m(M_m)\|_2)$$

$$= \langle \mathcal{P}_m(T(\mathcal{X})), \mathcal{P}_m(M_m) \rangle + O\left(\sqrt{kV/(N\tau)}\right) (\|\mathcal{P}_m(T(\mathcal{X}))\|_2 + \|M_m\|_2 + s) \qquad (\mathcal{R}_m \text{ is an isometry})$$

$$= \langle T(\mathcal{X}), \mathcal{P}_m(M_m) \rangle + O\left(\sqrt{kV/(N\tau)}\right) (\|\mathcal{P}_m(T(\mathcal{X}))\|_2 + \|M_m\|_2 + s) \qquad (\mathcal{P}_m \text{ is a projection})$$

$$= \langle T(\mathcal{X}), M_m \rangle + O\left(\sqrt{kV/(N\tau)}\right) (\|\mathcal{P}_m(T(\mathcal{X}))\|_2 + \|M_m\|_2 + s) + s \|\operatorname{proj}_{M_m - \mathcal{P}_m(M_m)}(T(\mathcal{X}))\|_2 .$$

Note that all of the computations before the last two steps of our algorithm are independent of $T(\mathcal{X})$, and thus $T(\mathcal{X})$ is independent of P_m and the vector $M_m - P_m(M_m)$. Since $P_m(T(\mathcal{X}))$ is the projection of $T(\mathcal{X})$ onto U_m , a k-dimensional subspace, we have that $\mathbf{E}_{\mathcal{X}\sim X}[\|P_m(T(\mathcal{X}))\|_2^2]^{1/2} = O(\sqrt{Vk})$. Since $M_m - P_m(M_m)$ is just a vector, the square root of the expected squared error of the projection of $T(\mathcal{X})$ onto that vector is $O(\sqrt{V})$.

By taking the expectation over $\mathcal{X} \sim X$ in the above, using Claim 3.7, (1) and (3) completes the proof of Proposition 1.7.

4 Learning Theory Applications

The structure of this section is as follows: We start in Section 4.1 with some technical machinery that is useful in all our subsequent learning applications. Section 4.2 presents our results on onehidden-layer neural networks. Section 4.3 presents our algorithmic results for mixtures of spherical Gaussians (both density and parameter estimation). Finally, Section 4.4 presents our algorithm for learning mixtures of linear regressions.

4.1 Technical Machinery

Before we get to the applications, we require the following technical tool that we briefly motivate. For our learning applications, we need to compute inner products of functions/distributions with Hermite polynomials. We have technology for computing inner products of tensors with objects that can be expressed as the expectation of sequential tensor computations. Unfortunately, the Hermite tensor $H_n(x)$ cannot conveniently be written in this form.

[LL22] resolves this issue by expressing $H_n(x)$ as a sum of tensors of x's and I's and replacing each copy of I by $y_i \otimes y_i$ for independent Gaussians y_i (noting that $\mathbf{E}[y_i \otimes y_i] = I$). We follow a different, more efficient approach: we take all of the copies of I in our product and replace them by a tensor power of a single random Gaussian vector y. This has the same expectation, but as a sum of only 2^n rank-1 terms — rather than roughly n^n . Furthermore, we can compute this quantity more cleverly using a sequential tensor computation of size only O(n).

Definition 4.1 (Extended Hermite Tensor). For $n \in \mathbb{N}$ and $x, y \in \mathbb{R}^d$, we define

$$H_n(x,y) = \frac{1}{\sqrt{n!}} \sum_{\substack{\text{Partitions of } [n] \text{ into } S_1, S_2\\ \text{with } |S_2| \text{ even}}} (-1)^{|S_2|/2} x^{\otimes S_1} \otimes y^{\otimes S_2} = \operatorname{Re}((x+iy)^{\otimes n})/\sqrt{n!} \, .$$

Note that $H_n(x, y)$ is an order-*n* tensor of dimension *d*. The important property of the above definition is that $H_n(x, y)$ can be computed by a sequential tensor computation of size O(n) and that $\mathbf{E}_{Y \sim N(0,I)}[H_n(x,Y)] = H_n(x)$. This follows from the fact that

$$\mathbf{E}[Y^{\otimes 2t}] = \sum_{\text{partitions } P \text{ of } [2t] \text{ into sets of size } 2} I^P .$$

Plugging this in to the expansion of $H_n(x, y)$ gives exactly the standard expansion of $H_n(x)$.

Remark 4.2. All examples in [LL22] involve explicitly writing the tensors as sums of rank-1 tensors (i.e., products of vectors). While these are all STCs, not all STCs can be efficiently written in this way. For example, if expanded out, $H_n(x, y)$ can be written as a sum of 2^n rank-1 tensors. However, thinking of it as the real part of $(x + iy)^{\otimes n}$, we can use a dynamic program keeping track of the real and imaginary parts, and write it as the output of a STC of size O(n) instead.

To apply Proposition 1.7, we also need to bound the second moment of $H_n(X, Y)$ over $Y \sim N(0, I)$ and independent $X \sim N(\mu, I)$.

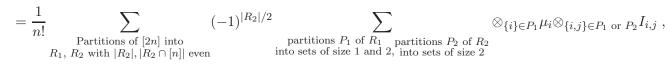
Lemma 4.3 (Second Moment Bound of $H_n(X, Y)$). We have that $\mathbf{E}_{Y,X}[H_n(X, Y) \otimes H_n(X, Y)]$, where $X \sim N(\mu, I)$, $Y \sim N(0, I)$ and X, Y are independent, is bounded above by $O(||\mu||_2^2/n + 1)^n$.

Proof. The expectation of $H_n(X,Y) \otimes H_n(X,Y)$ can be expressed as follows:

$$\begin{split} \mathbf{E}_{X,Y}[H_n(X,Y) \otimes H_n(X,Y)] &= \\ & \frac{1}{n!} \sum_{S_1,S_2 \text{partition of } [n]} (-1)^{|S_2|/2} x^{\otimes S_1} y^{\otimes S_2} \otimes \sum_{T_1,T_2 \text{ partition of } [n]} (-1)^{|T_2|/2} x^{\otimes T_1} y^{\otimes T_2} \;. \end{split}$$

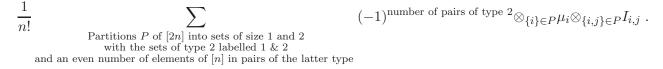
Viewing the above as a partition of [2n] into $R_1 = S_1 \cup T_1$ and $R_2 = S_2 \cup T_2$, we can write

$$\mathbf{E}_{X,Y}[H_n(X,Y) \otimes H_n(X,Y)] = \frac{1}{n!} \sum_{\substack{\text{Partitions of } [2n] \text{ into}\\R_1, R_2 \text{ with } |R_2|, |R_2 \cap [n]| \text{ even}}} (-1)^{|R_2|/2} x^{\otimes R_1} y^{\otimes R_2}$$



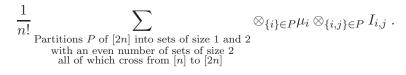
where the μ_i denotes a μ in the *i*-th tensor slot and $I_{i,j}$ denotes a copy of the identity in the *i*th and *j*th.

If we now consider $P_1 \cup P_2$, we get a partition of [2n] into sets of size 1 and 2 with the sets of size 2 coming from either P_1 or P_2 (which we label as type 1 and type 2). We have the restriction that the number of elements of [n] in pairs of type 2 is even. Thus, we get



We think of this as first choosing the partition P of [2n] into sets of size 1 and 2 and only later choosing which to label as type 1 and type 2. We note that if there is any pair contained entirely in [n] or entirely in $[2n] \setminus [n]$, then switching its type from 1 to 2 or back keeps the term the same but reverses its sign. Thus, these terms cancel out.

Hence, we have:



Considering terms with exactly k sets of size 2, there are $\binom{n}{k}^2$ ways to choose the sets of size 1, and k! ways to choose how to pair up the remaining elements; and at most 2^k ways to label the sets of size 2. When thought of as a $d^n \times d^n$ matrix, the tensor in the sum then has spectral norm at most $\|\mu\|_2^{2n-2k}$. Thus, the spectral norm of the covariance is at most

$$\frac{2^{O(n)}}{n^n} \sum_k \|\mu\|_2^{2n-2k} \sqrt{n^{2k}} = O(\|\mu\|_2^2/n+1)^n ,$$

completing the proof.

Remark 4.4. For the case of a mixture of k identity covariance Gaussians, the above second moment bound is polynomial in k times the largest mean, even when $n = O(\log(k))$. This bound is the key ingredient that allows us to obtain the optimal mean separation condition of $O(\sqrt{\log(k)})$ for the parameter estimation problem, improving on the bound of [LL22].

4.2 Learning Positive Linear Combinations of ReLU Activations

Problem Setup We start by formally defining the target class of functions to be learned.

Definition 4.5 (Positive Linear Combinations of ReLUs). Let $C_{d,k}$ denote the concept class of onehidden-layer ReLU networks on \mathbb{R}^d with k hidden units and positive coefficients. That is, $F \in C_{d,k}$ if and only if there exist k unit vectors $v_i \in \mathbb{R}^d$ and non-negative coefficients $w_i \in \mathbb{R}_+$, $i \in [k]$, with $\sum_{i=1}^k w_i = 1$ such that $F(x) = \sum_{i=1}^k w_i \text{ReLU}(v_i \cdot x)$, where $\text{ReLU}(t) = \max\{0, t\}, t \in \mathbb{R}$.

Following prior work on this problem [DKKZ20, GGJ⁺20, DK20, CKM21, CGKM22, CDG⁺23, DK24], we will assume that the feature vectors are normally distributed. The definition of the PAC learning problem in our setting is the following.

Definition 4.6 (PAC Learning $C_{d,k}$). The PAC learning problem for the class $C_{d,k}$ is the following: The input is a multiset of i.i.d. labeled examples (x, y), where $x \sim N(0, I)$ and y = F(x), for an unknown $F \in C_{d,k}$. The goal of the learner is to output a hypothesis $H : \mathbb{R}^d \to \mathbb{R}$ that with high probability is close to F in L_2 -norm, i.e., satisfies $||H - F||_2 \leq \epsilon ||F||_2$. The hypothesis H is allowed to lie in any efficiently representable hypothesis class \mathcal{H} .

Our main algorithmic result in this context is the following.

Theorem 4.7 (Learning Algorithm for $C_{d,k}$). There is a PAC learning algorithm for $C_{d,k}$ with respect to the standard Gaussian distribution on \mathbb{R}^d with the following performance guarantee: Given $\epsilon > 0$, and access to labeled examples from an unknown target $F \in C_{d,k}$, the algorithm has sample and computational complexity $\operatorname{poly}(d,k)2^{\operatorname{poly}(1/\epsilon)}$, and outputs a hypothesis $H : \mathbb{R}^d \to \mathbb{R}$ that with high probability satisfies $||H - F||_2 \le \epsilon ||F||_2$.

It is easy to see that our results straightforwardly extend to the case that the labels have been corrupted by random zero-mean additive noise.

As will become clear from the analysis of this subsection, our PAC learning algorithm is not proper. The hypothesis H output by our algorithm is a succinct description of a low-degree polynomial (namely, of degree poly $(1/\epsilon)$). Specifically, each Hermite coefficient will be expressed in compressed form, as per Proposition 1.7. This succinct description allows us to produce an efficient evaluation oracle for H(x).

To apply Proposition 1.7 for this learning problem, we will need the following basic facts on the Hermite decomposition of the relevant functions.

Hermite Decomposition of Functions in $C_{d,k}$ We will require the following lemma:

Lemma 4.8 (Hermite Expansion of $C_{d,k}$). Let $F : \mathbb{R}^d \to \mathbb{R}$ be any function of the form $F(x) = \sum_{i=1}^k w_i \operatorname{ReLU}(v_i \cdot x)$ for some $v_i \in \mathbb{R}^d$ with $||v_i||_2 = 1$ and $w_i \in \mathbb{R}_+$. Then F has Hermite expansion $F(x) = \sum_{n=0}^{\infty} \langle T_n, H_n(x) \rangle$, where H_n is the normalized Hermite tensor (Definition 2.2) and $T_n \in (\mathbb{R}^d)^{\otimes n}$ is defined by

$$T_n \stackrel{\text{def}}{=} c_n \sum_{i=1}^k w_i v_i^{\otimes n} , \qquad (4)$$

where c_n are explicit constants defined as follows: If n > 1, then $c_n = 0$ if n is odd and

$$c_n = (-1/4)^{(n-2)/4} \sqrt{\binom{n-2}{(n-2)/2}} / \sqrt{2\pi n(n-1)} = \Theta(n^{-5/4})$$

if n is even.

Proof. These calculations have essentially appeared in prior work. For the sake of completeness, here we show how to directly deduce the lemma from prior work.

By orthogonality of the Hermite tensors, we have that the tensor T_n in the Hermite expansion of F, $F(x) = \sum_{n=0}^{\infty} \langle T_n, H_n(x) \rangle$, is defined by $T_n = \mathbf{E}_{X \sim N(0,I)} [F(X)H_n(X)]$.

Lemma 3.1 of [DK24] shows that for all $n \in \mathbb{Z}_+$, we have that $\mathbf{E}_{G \sim N(0,1)}[\operatorname{ReLU}(G)h_n(G)] = c_n$ for the values of c_n in the lemma statement. This in turn implies that for any unit vector $v \in \mathbb{R}^d$, it holds that $\operatorname{ReLU}(v \cdot x) = \sum_{n=0}^{\infty} c_n \langle H_n(x), v^{\otimes n} \rangle$ (Lemma 3.2 of [DK24]). Via orthogonality, we additionally get that $\mathbf{E}_{X \sim N(0,I)}[\operatorname{ReLU}(v \cdot X)H_n(X)] = c_n v^{\otimes n}$ (Corollary 3.3 of [DK24]). Since $F(x) = \sum_{i=1}^{k} w_i \text{ReLU}(v_i \cdot x)$, by linearity of expectation it follows that

$$\mathbf{E}_{X \sim N(0,I)}[F(X)H_n(X)] = \sum_{i=1}^k w_i \mathbf{E}_{X \sim N(0,I)}[\text{ReLU}(v_i \cdot X)H_n(X)] = \sum_{i=1}^k w_i c_n v_i^{\otimes n} = c_n \sum_{i=1}^k w_i v_i^{\otimes n} .$$

This completes the proof.

Proof of Theorem 4.7 The key properties that enable our algorithm are that the tensor $H_n(x, y)$ of Definition 4.1 (i) can be computed by a sequential tensor computation of size O(n), (ii) it satisfies $\mathbf{E}_{Y \sim N(0,I)}[H_n(x,Y)] = H_n(x)$, for any $x \in \mathbb{R}^d$, and (iii) it has bounded second moment (Lemma 4.3).

To leverage these properties in the context of our learning problem, we note that the second property implies that

$$T_n = \mathbf{E}_{X \sim N(0,I)}[F(X)H_n(X)] = \mathbf{E}_{X \sim N(0,I),Y \sim N(0,I)}[F(X)H_n(X,Y)],$$
(5)

where X and Y are independent. Indeed, the first equality is the definition of T_n and the second follows using property (ii) above (by linearity of expectation).

Lemma 4.8 (namely, the definition of the c_n 's) implies that T_n is zero, if n > 1 and n odd. Thus, it suffices to consider the Hermite coefficients of order n = 1 or n even. We claim that, for n = 1or n even, we can apply Proposition 1.7 to efficiently approximate T_n . Specifically, we will show how to efficiently approximate the inner product $\langle T_n, H_n(x) \rangle$ for $x \sim N(0, I)$. Given that such an efficient computation is possible, the hypothesis of our learning algorithm will be an approximation to the low-degree Hermite expansion of F of the function

$$\widetilde{H}(x) = \sum_{n=0}^{n_0} \langle \widetilde{T_n, H_n(x)} \rangle , \qquad (6)$$

where $n_0 = \text{poly}(1/\epsilon)$. Note that if $H(x) = \sum_{n=0}^{n_0} \langle T_n, H_n(x) \rangle$ is the (exact) low-degree Hermite expansion of F up to degree n_0 , them $\mathbf{E}_{x \sim N(0,I)}[(H(x) - F(x))^2] \leq (\epsilon/2)^2$. So for our purposes it suffices to approximate each term in the sum up to error $\epsilon/(2n_0) = \text{poly}(\epsilon)$.

We will show how to obtain an evaluation oracle for H with sample and computational complexity $\operatorname{poly}(d,k)2^{\operatorname{poly}(1/\epsilon)}$. To do this, we will apply Proposition 1.7 to obtain an approximation $\langle T_n, H_n(x) \rangle$ of $\langle T_n, H_n(x) \rangle$ for all $n \leq n_0$ (excluding odd values greater than one). It turns out that the sample size and running time are dominated by the case $n = n_0$.

We proceed to apply Proposition 1.7. The parameter m in the statement of the proposition will be set to m := n, for $n \le n_0$. The tensor M_t , for $t \le 2m$ with t even or equal to m, will be set to $M_t := \sum_{i=1}^t w_i v_i^{\otimes t}$, where w_i and v_i are the hidden weights and parameter vectors of the target function F(x).

We exhibit a sequential tensor computation S_t with the desired properties. Specifically, S_t is a sequential tensor computation with three inputs: two vector inputs $v, u \in \mathbb{R}^d$ and a scalar input $y \in \mathbb{R}$. The output of $S_t(u, v, y)$ is the order-*t* and dimension-*d* tensor $(1/c_t)yH_t(u, v)$, where $c_t = \Theta(t^{-5/4})$ is the explicit constant in Lemma 4.8. Since $H_t(u, v)$ can be computed by a sequential tensor computation of size O(t), the same holds for $(1/c_t)yH_t(u, v)$.

The input distribution corresponding to these inputs is the joint distribution $\mathcal{D} := (X, Y, F(X))$, where X and Y are independent standard Gaussians, $X, Y \sim N(0, I)$, and F(X) is the value of the target function. Since we have sample access to labeled examples from the target function, i.e., from the distribution (X, F(X)), the distribution \mathcal{D} is efficiently samplable. We now have that

$$\mathbf{E}_{(u,v,y)\sim\mathcal{D}}\mathcal{S}_t(u,v,y) = \mathbf{E}_{X\sim N(0,I),Y\sim N(0,I)}[(1/c_t)F(X)H_t(X,Y)] = (1/c_t)T_t = \sum_{i=1}^k w_i v_i^{\otimes t} ,$$

where the first equality holds by definition, the second equality follows from (5), and the third equality follows from Lemma 4.8 (Equation (4)).

This gives us a sequential tensor computation S_t of size O(t) with expectation $M_t = \sum_{i=1}^k w_i v_i^{\otimes t}$, when t is equal to 1 or even. It remains to bound above the covariance of S_t . We show the following:

Claim 4.9. We have that $\operatorname{Cov}_{(u,v,y)\sim\mathcal{D}}[\mathcal{S}_t(u,v,y)]$ is bounded above by $V := O(1)^t$.

Proof. Recall that for $(u, v, y) \sim \mathcal{D}$, we have that $\mathcal{S}_t(u, v, y) = (1/c_t)F(X)H_t(X, Y)$. By definition, we want to bound from above the variance of the random variable $\langle W, (1/c_t)F(X)H_t(X, Y)\rangle = (1/c_t)F(X)\langle W, H_t(X, Y)\rangle$ for any tensor $W \in (\mathbb{R}^d)^{\otimes t}$ with $||W||_2 = 1$. We can write

$$\begin{aligned} \mathbf{Var}[\langle W, (1/c_t)F(X)H_t(X,Y)\rangle] &= (1/c_t)^2 \mathbf{Var}[\langle W,F(X)H_t(X,Y)\rangle] \\ &= \Theta(t^{5/2})\mathbf{Var}[F(x)\langle W,H_t(X,Y)\rangle] \\ &\leq \Theta(t^{5/2}) \|F\|_4^2 \|\langle W,H_t(X,Y)\rangle\|_4^2 \\ &\leq \Theta(t^{5/2}) O(1) O(1)^t \|\langle W,H_t(X,Y)\rangle\|_2^2 \\ &\leq \Theta(t^{5/2}) O(1) O(1)^t O(1)^t \\ &= O(1)^t , \end{aligned}$$

where the second line follows from the definition of c_t (for t = 1 or even), and the third line uses Cauchy-Schwartz. The fourth line follows from two facts. First, for any $m \in \mathbb{Z}_+$, we have

$$\|F(X)\|_m \le \sum_{i=1}^k w_i \|\operatorname{ReLU}(v_i \cdot X)\|_m \le \sum_{i=1}^k w_i \|v_i \cdot X\|_m = O(\sqrt{m}) \sum_{i=1}^k w_i = O(\sqrt{m}) .$$

Second, since $\langle W, H_t(X, Y) \rangle$ is a degree-*t* polynomial over Gaussians, hypercontractivity (Fact 2.4) implies that this quantity is $O(1)^t || \langle W, H_t(X, Y) \rangle ||_2^2$.

Finally, by Lemma 4.3 (applied for $\mu = 0$), it follows that $\|\langle W, H_t(X, Y) \rangle\|_2^2$ is at most $O(1)^t$. This justifies the fifth line and completes the proof of the claim.

It remains to exhibit the sequential tensor computation \mathcal{F}_m , where we recall that m = n. \mathcal{F}_n is a sequential tensor computation with two input vector inputs $v, u \in \mathbb{R}^d$. The output of $\mathcal{F}_n(u, v)$ is the order-*n* and dimension-*d* tensor $H_n(u, v)$. We recall the fact that $H_n(u, v)$ can be computed by a sequential tensor computation of size O(n).

The input distribution corresponding to these inputs is the (X, Y), where X and Y are independent standard Gaussians, $X, Y \sim N(0, I)$. We note that \mathcal{D}' is the distribution of Y (that we have sample access to) while X is a single Gaussian sample (see the proposition statement). Since $H_n(X) = \mathbf{E}_Y[H_n(X,Y)]$ and the covariance of $H_n(X,Y)$ is bounded above by $O(1)^n$ (by Lemma 4.3), we have all the ingredients we need to apply Proposition 1.7.

By Proposition 1.7, if the total number of samples from \mathcal{D} and \mathcal{D}' is $N = \text{poly}(d, k, S, V, 1/\epsilon)$, where $V = O(1)^n$, for a Gaussian random X we can approximate $\langle T_n, H_n(X) \rangle$ up to expected squared error ϵ in poly(N, d) time.

Recall that $V = O(1)^n$ and S = O(n). By setting $n = n_0 = \text{poly}(1/\epsilon)$ and the desired accuracy ϵ to be $\epsilon/n_0 = \text{poly}(1/\epsilon)$, we obtain estimates that suffice to evaluate the hypothesis \widetilde{H} .

Overall, this gives a learning algorithm with sample complexity and runtime $poly(d, k)2^{poly(1/\epsilon)}$, completing the proof of Theorem 4.7.

4.3 Learning Mixtures of Spherical Gaussians

In Section 4.3.1, we give our density estimation algorithm. In Section 4.3.2, we give our parameter estimation algorithm.

Problem Setup We start by formally defining the target distribution class to be learned.

Definition 4.10 (Mixtures of Spherical Gaussians). A k-mixture of spherical Gaussians is a distribution on \mathbb{R}^d defined by $F = \sum_{i=1}^k w_i N(\mu_i, I)$, where $\mu_i \in \mathbb{R}^d$ are the unknown mean vectors and $w_i \ge 0$, with $\sum_{i=1}^k w_i = 1$, are the mixing weights.

We will consider both density estimation and parameter estimation. In density estimation, we want to output a hypothesis distribution whose total variation distance from the target distribution is small. Recall that an ϵ -sampler for a distribution D is a circuit C that on input a set z of uniformly random bits it generates then $y \sim D'$, for some distribution D' which has $d_{\text{TV}}(D', D) \leq \epsilon$.

Definition 4.11 (Density Estimation for Mixtures of Spherical Gaussians). The density estimation problem for mixtures of spherical Gaussians is the following: The input is a multiset of i.i.d. samples in \mathbb{R}^d drawn from an unknown k-mixture $F = \sum_{i=1}^k w_i N(\mu_i, I)$. The goal of the learner is to output a (sampler for a) hypothesis distribution H such that with high probability $d_{\text{TV}}(H, F) \leq \epsilon$.

Definition 4.12 (Parameter Estimation for Mixtures of Spherical Gaussians). The parameter estimation problem for mixtures of spherical Gaussians is the following: The input is a multiset of i.i.d. samples in \mathbb{R}^d drawn from an unknown k-mixture $F = \sum_{i=1}^k w_i N(\mu_i, I)$, where the component means μ_i satisfy a pairwise separation condition. The goal of the learner is to accurately estimate the weights and mean vectors of the components.

4.3.1 Density Estimation for Mixtures of Spherical Gaussians

We consider mixtures of k identity covariance Gaussians on \mathbb{R}^d with the additional restriction that the mean vectors of the components lie in a ball of radius $O(\sqrt{\log(k)})$. By subtracting the mean of the mixture, it suffices to consider the case that each component mean has magnitude $O(\sqrt{\log(k)})$. Our main algorithmic result in this context is the following.

Theorem 4.13 (Density Estimation Algorithm for Mixtures of Spherical Gaussians with Bounded Means). There is an algorithm that given $\epsilon > 0$ and $n = \text{poly}(d, k, 1/\epsilon)$ samples from an unknown $F = \sum_{i=1}^{k} w_i N(\mu_i, I)$ on \mathbb{R}^d with means of magnitude $O(\sqrt{\log(k)})$, it runs in poly(n, d) time and outputs a (sampler for a) hypothesis distribution H such that with high probability $d_{\text{TV}}(H, F) \leq \epsilon$.

As for our previous application, our learning algorithm is not proper. The hypothesis distribution H will be such that H/G, where G is the pdf of N(0, I), will be a low-degree polynomial in compressed form, as per Proposition 1.7, that allows for efficient sampling.

To apply Proposition 1.7 for this learning problem, we will need the following basic facts on the Hermite decomposition of the relevant distributions.

Hermite Decomposition of Spherical Mixtures We will require the following lemma:

Lemma 4.14 (Hermite Expansion of Spherical Mixtures). Let $F : \mathbb{R}^d \to \mathbb{R}$ be any distribution of the form $F(x) = \sum_{i=1}^k w_i N(\mu_i, I)$ for some $\mu_i \in \mathbb{R}^d$ and $w_i \in \mathbb{R}_+$ with $\sum_{i=1}^k w_i = 1$. Let G(x)be the pdf of the standard Gaussian N(0, I). Then (F/G)(x) has Hermite expansion (F/G)(x) = $\sum_{n=0}^{\infty} \langle T_n, H_n(x) \rangle$, where H_n is the normalized Hermite tensor (Definition 2.2) and $T_n \in (\mathbb{R}^d)^{\otimes n}$ is defined by

$$T_n \stackrel{\text{def}}{=} (1/\sqrt{n!}) \sum_{i=1}^k w_i \mu_i^{\otimes n} .$$
(7)

Proof. Since (F/G) is in L_2 with respect to the Gaussian measure G, a Hermite expansion does exist in the form of $(F/G)(x) = \sum_{n=0}^{\infty} \langle T_n, H_n(x) \rangle$, where

$$T_n = \mathbf{E}_{x \sim G}[(F/G)(x)H_n(x)] = \mathbf{E}_{x \sim F}[H_n(x)].$$

Since F is a mixture, we have that $\mathbf{E}[H_n(F)] = \sum_i w_i \mathbf{E}[H_n(N(\mu_i, I))]$. By Lemma 2.7 of [Kan21], this is equal to $(1/\sqrt{n!}) \sum w_i \mu_i^{\otimes n}$. This completes the proof.

Proof of Theorem 4.13 The proof is analogous to the proof of Theorem 4.7 for our previous application. We similarly leverage the fact that the tensor $H_n(x, y)$ of Definition 4.1 (i) can be computed by a sequential tensor computation of size O(n), (ii) it satisfies $\mathbf{E}_{Y \sim N(0,I)}[H_n(x,Y)] = H_n(x)$, for any $x \in \mathbb{R}^d$, and (iii) it has bounded second moment (Lemma 4.3).

We start by noting that the second property implies that the tensor T_n of (7) satisfies

$$T_n = \mathbf{E}_{X \sim F}[H_n(X)] = \mathbf{E}_{X \sim F, Y \sim N(0,I)}[H_n(X,Y)], \qquad (8)$$

where X and Y are independent. Indeed, the first equality is the definition of T_n and the second follows using property (ii) above (by linearity of expectation).

We will apply Proposition 1.7 to efficiently approximate T_n . Specifically, we will show how to efficiently approximate the inner product $\langle T_n, H_n(x) \rangle$ for $x \sim N(0, I)$. Given that such an efficient approximation is possible, we will show how to output an efficiently samplable hypothesis distribution.

The first step will be to produce an evaluation circuit which computes an approximation to (F/G)(x) for Gaussian random x. We will achieve this via an application of Proposition 1.7.

Note that

$$||T_n||_2^2 = O\left(\max_{i \in [k]} ||\mu_i||_2^2/n\right)^{n/2} = O(\log(k)/n)^{n/2},$$

where the first equality follows by a direct calculation and the second follows by our assumed bound on the $\|\mu_i\|_2$'s.

Therefore, for n_0 a sufficiently large constant multiple of $\log(k/\epsilon)$, truncating the Hermite expansion of F/G at n_0 introduces L_2 -error (with respect to G) at most $\epsilon/2$, i.e.,

$$(\widetilde{F}/G)(x) = \sum_{n=0}^{n_0} \langle T_n, H_n(x) \rangle$$

is $\epsilon/2$ -close to F/G in L_2 norm. In particular, if we let $\widetilde{F}(x) = G(x) \sum_{n=0}^{n_0} \langle T_n, H_n(x) \rangle$, we have that the L_1 -error between F and \widetilde{F} is less than $\epsilon/2$.

Using Proposition 1.7, we can come up with an evaluation circuit which approximates (\bar{F}/G) at random x to small error. The parameter m in the statement of the proposition will be set to m := n, for $n \leq n_0$. The tensor M_t will be set to $M_t := \sum_{i=1}^k w_i v_i^{\otimes t}$, where w_i are the mixture weights and $v_i = \mu_i / \sqrt{\log(k)}$, where μ_i are the component means of F(x). By our assumption on

the magnitude of the μ_i 's, it follows that each v_i has ℓ_2 -norm O(1). This re-parameterization is necessary in order to get the right error terms in our analysis. We then have that

$$(\widetilde{F}/G)(x) = \sum_{n=0}^{n_0} \langle \sqrt{\log(k)^n/n!} H_n(x), M_n \rangle .$$

Combining the above with (8) gives that $M_n = \sqrt{n!/\log(k)^n} \mathbf{E}[H_n(F,G)].$

Given the above, the sequential tensor computation S_n has two vector inputs $v, u \in \mathbb{R}^d$ and is defined so that $S_n(u,v) = \sqrt{n!/\log(k)^n}H_n(u,v)$. We have already argued that $H_n(u,v)$ can be computed by a sequential tensor computation of size O(n). The corresponding input distribution is the joint distribution $\mathcal{D} := (X,Y)$, where $X \sim F$ and $Y \sim N(0,I)$ are independent. Since we have sample access to F the distribution \mathcal{D} is efficiently samplable. Moreover, we have

$$\mathbf{E}_{(u,v,)\sim\mathcal{D}}[\mathcal{S}_n(u,v)] = \sqrt{n!/\log(k)^n} \mathbf{E}_{X\sim F, Y\sim N(0,I)}[H_n(X,Y)] = M_n \; .$$

We can use Lemma 4.3 to bound the covariance of S_t . We show the following:

Claim 4.15. We have that $\mathbf{Cov}_{(u,v)\sim\mathcal{D}}[\mathcal{S}_t(u,v)]$ is bounded above by $V := \mathrm{poly}(k/\epsilon)$.

Proof. Using the fact that $F = \sum_{i=1}^{k} w_i N(\mu_i, I)$, we find that the covariance of $H_n(F, G)$ is $\sum_{i=1}^{k} w_i \mathbf{Cov}[H_n(N(\mu_i, I), G)] + \mathbf{Cov}[X]$, where X is the random variable that has expected value $\mathbf{E}[H_n(N(\mu_i, I), G)] = \mu_i^{\otimes n}/\sqrt{n!}$ with probability w_i . By Lemma 4.3, the first term above has operator norm at most $O(\log(k)/n+1)^n$ and the second term is similarly bounded. For $n = O(\log(k/\epsilon))$, this is at most $\operatorname{poly}(k/\epsilon)$.

The sequential tensor computation \mathcal{F}_n and the distribution \mathcal{D}' will be the same as in our previous application, up to rescaling. Namely, $\mathcal{F}_n(u,v) = \sqrt{\log(k)^n/n!}H_n(u,v)$ and \mathcal{D}' is the distribution of $Y \sim N(0,I)$. Since $H_n(X) = \mathbf{E}_Y[H_n(X,Y)]$, for $X \sim N(0,I)$, and the second moment of $H_n(X,Y)$ is bounded above by Lemma 4.3, it follows that the second moment of \mathcal{F}_n is bounded above by $O(\log(k)/n)^n$.

Therefore, applying Proposition 1.7, we have an algorithm that runs in time $\operatorname{poly}(k, d, n_0, 1/\epsilon, 1/\tau)$ which with probability $1 - \tau$ produces an approximation to $\langle \sqrt{\log(k)^n/n!}H_n(x), M_n \rangle$, for $n \leq n_0$, so that the L_2 (and thus L_1) expected error for $x \sim G$ is at most $\epsilon/(2(n_0 + 1))$. Summing this over $0 \leq n \leq n_0$ and taking $\tau < 1/(100n_0)$, gives a 99% probability of producing an approximation to (\tilde{F}/G) with L_1 error at most $\epsilon/2$. This in turn gives an approximation to F/G with L_1 error at most ϵ .

We next need to address how to go from an approximate evaluation oracle for F/G to an approximate sampler for F. We achieve this via rejection sampling. In particular, we produce a random sample $x \sim G$ and compute our oracle R at x. We then return x with probability proportional to R(x), and otherwise resample and repeat. Note that if R returned exactly (F/G) and if we could keep x with probability exactly proportional to R(x), this would give F exactly.

The first issue we have to face is that R(x) is neither guaranteed to be non-negative nor bounded, and thus keeping x with probability proportional to R(x) is impossible. To fix this, we need to truncate R. In particular, we define R' so that R' = 0 if R < 0; R = A if R > A, and R otherwise, for A some parameter we will choose shortly.

We claim that $||R' - (F/G)||_{1,G} := \mathbf{E}_{x \sim G}[|R(x) - (F/G)(x)|]$ is at most 3ϵ .

First, note that $||R - (F/G)||_{1,G} < \epsilon$. Next, since (F/G) is non-negative, replacing R(x) by 0 when R(x) < 0, will only decrease the distance. It remains to show that

$$\mathbf{E}_{x\sim G}[\max(R(x) - A, 0)] < 2\epsilon \; .$$

Since $||R - (F/G)||_{1,G} < \epsilon$, it suffices to show that $\mathbf{E}_{x \sim G}[\max((F/G)(x) - A, 0)] < \epsilon$. The above is equal to $\int_{t>A} \mathbf{Pr}[(F/G)(x) > t] dt$. Using the definition of F, it follows that

$$(F/G)(x) \le \max_{i} \frac{N(\mu_i, I)}{N(0, I)}(x) \le \max_{i} \exp(\mu_i \cdot x) .$$

Thus, by a union bound

$$\mathbf{Pr}[(F/G)(x) > t] \le \sum_{i} \mathbf{Pr}[\mu_i \cdot x > \log(t)] \le O(k) \exp\left(-\Omega\left(\log(t)/\sqrt{\log(k)}\right)^2\right)$$

Thus, if A is a large constant degree polynomial in k/ϵ , this resulting integral is at most ϵ .

If the above holds, then our algorithm can evaluate a function R'(x) so that with 99% probability $||R - (F/G)||_1 < 3\epsilon$ (we can re-parameterize, by setting $\epsilon/6$ instead of ϵ , so this is actually less than $\epsilon/2$); and $R'(x) \in [0, A]$ for all x.

We now use rejection sampling. We repeatedly sample $x \sim G$ and with probability R'(x)/A, we return x and otherwise repeat this process. Note that since $||R' - F/G||_1$ is small, $\mathbf{E}[R'(x)] > 1/2$, and so each iteration has an $\Omega(1/A)$ chance of terminating. Thus, in expectation, this procedure terminates after only O(A) rounds with 99% probability. Furthermore, the resulting distribution that we sample from is proportional to R'G. Observe that this is ϵ -close to (F/G)G = F. Thus, we are sampling from a distribution ϵ -close to F, as desired.

Overall, this gives a density estimation algorithm with complexity $poly(d, k, 1/\epsilon)$, completing the proof of Theorem 4.13.

4.3.2 Parameter Estimation for Mixtures of $O(\sqrt{\log(k)})$ -Separated Gaussians

We start by recalling that [LL22]'s result on clustering/parameter estimation for mixtures of spherical Gaussians only works under pairwise mean separation of $\log(k)^{1/2+c}$, for some constant $c > 0^4$. This is because, due to the way that they were approximating Hermite polynomials, they could only approximate tensors of order $\log(k)/\log\log(k)$ in polynomial time. Our extended Hermite tensors (Definition 4.1 and Lemma 4.3) improve upon this, allowing us to work with tensors of order $\log(k)$, and thus obtain optimal (up to constant factors) $O(\sqrt{\log(k)})$ separation.

Specifically, we show the following statement that also handles general weights:

Theorem 4.16 (Parameter Estimation for Mixtures of Spherical Gaussians). Let $F = \sum_{i=1}^{k} w_i N(\mu_i, I)$ be a mixture of Gaussians with $w_{\min} \leq \min w_i$ and let $\alpha > 0$ be an accuracy parameter. Suppose that $s := \min_{i \neq j} ||\mu_i - \mu_j||_2$ is at least a sufficiently large constant multiple of $\sqrt{\log(1/(\alpha w_{\min}))}$ and furthermore that $\max_{i \neq j} ||\mu_i - \mu_j||_2 = O(\min_{i \neq j} ||\mu_i - \mu_j||_2)$. Then there exists an algorithm that given k, w_{\min}, α, s and $\operatorname{poly}(d/(w_{\min}\alpha))$ i.i.d. samples from F, runs in $\operatorname{poly}(n, d)$ time, and outputs estimates $\tilde{\mu}_i$ and \tilde{w}_i of μ_i and w_i , so that with probability 2/3, for some permutation π of [k]

$$|\tilde{w}_i - w_i|, \|\tilde{\mu}_i - \mu_i\|_2 \le \alpha$$

for all $i \in [k]$.

Note that [LL22] achieves this goal only for separation $\log(1/(\alpha w_{\min}))^{1/2+c}$. While they do not require the condition that the largest pairwise distance is comparable to the smallest, in order to remove this condition, they developed and leveraged a complicated recursive clustering argument (see Sections 10 and 11 of [LL22]), which we expect can also be applied to our setting.

⁴For the sake of simplicity, this description focuses on the case of uniform mixtures.

Proof of Theorem 4.16. Firstly, we note that by standard techniques, we can reduce to the case of d = k; see, e.g., [VW02, DKS18]. In particular, we may assume that $\max ||\mu_i - \mu_j||_2 < \log(k)$ (or otherwise [LL22]'s result will apply anyway). Given this and $poly(d/(w_{\min}\alpha))$ samples, we can approximate the covariance matrix of F to inverse polynomial accuracy. Letting V be the span of the k principal eigenvectors of this covariance matrix, it is not hard to show that all of the μ_i lie within distance $\alpha/2$ of some translate of V (this translate can be estimated by approximating the mean of the projection of F onto V^{\perp}). Thus, it suffices to consider our algorithm on the projection of F onto V, which is a k-dimensional subspace.

We next draw $poly(k/(w_{min}\alpha))$ samples from F. We claim that with high probability we can cluster these samples so that the samples drawn from the component $N(\mu_i, I)$ lie exactly in their own cluster. If we can do this, it is straightforward to show that letting $\tilde{\mu}_i$ be the mean of the elements in the i^{th} cluster and letting \tilde{w}_i be the fraction of elements in the i^{th} cluster will suffice for our estimates. Following [LL22], in order to do this clustering, we need only the following procedure:

Given $x, x' \sim F$ and $\operatorname{poly}(1/(w_{\min}\tau))$ i.i.d. samples from F (for some $\tau = \operatorname{poly}(k/(w_{\min}\alpha)) > 0$), determine with probability at most $1 - \tau$ of error whether or not x and x' were drawn from the same component.

For this, we consider the distribution X' obtained by taking $(y - y')/\sqrt{2}$, where y, y' are independent samples from F. We note that X' is a mixture of at most k^2 Gaussians, namely

$$X' \sim (w_1^2 + w_2^2 + \ldots + w_k^2)N(0, I) + \sum_{i \neq j} w_i w_j N((\mu_i - \mu_j)/\sqrt{2}, I)$$

We define $M_t = \sum_{i,j} w_i w_j (\mu_i - \mu_j)^{\otimes t}$ and note that

$$M_t = \sqrt{t!} \mathbf{E}_{x \sim X', y \sim N(0, I)} [H_t(x, y)] .$$
(9)

Our goal will be to use Proposition 1.7 to approximate $P_t := \langle H_t((x - x')/\sqrt{2}), M_t \rangle$, for t an even integer bigger than a suitably large constant multiple of $\log(1/w_{\min})$. In particular, we have that

$$P_t = \sum_{i \neq j} w_i w_j \|\mu_i - \mu_j\|_2^t h_t(v_{i,j} \cdot (x - x')/\sqrt{2}) ,$$

where $v_{i,j}$ is the unit vector in the direction of $\mu_i - \mu_j$.

If x comes from the i^{th} component and x' from the j^{th} , $v_{i,j} \cdot (x - x')/\sqrt{2}$ is distributed as $N(\|\mu_i - \mu_j\|_2/\sqrt{2}, 1)$, and thus with high probability is at least s/2. By direct computation, we have that

$$h_t(x) = \frac{1}{\sqrt{t!}} \sum_{a=0}^{t/2} (-1)^a x^{t-2a} \binom{t}{2, 2, 2, \dots, 2, t-2a} = \frac{x^t}{\sqrt{t!}} (1 + O((t/x)^2 + (t/x)^4 + \dots + (t/x)^t)).$$

Therefore, if $|x| \leq t$, we have that $h_t(x) = O(t)^{t/2}$; but if |x| is at least a sufficiently large constant multiple of t, we have that $h_t(x) \geq \frac{x^t}{2\sqrt{t!}}$. In particular, $h_t(v_{i,j} \cdot (x - x')/\sqrt{2})$ will be at least $\Omega(s^2/\sqrt{t})^t$. Since all of the other $h_t(v_{i',j'} \cdot (x - x')/\sqrt{2})$ are bounded below, this implies that $P_t = w_{\min}^2 \Omega(s^2/\sqrt{t})^t$.

On the other hand if x and x' come from the same component, $v_{i,j}(x-x')/\sqrt{2} \sim N(0,1)$, and with probability $1-\tau/3$ all of these quantities have magnitude $O(\sqrt{\log(k/\tau)})$. Direct computation similarly gives us that

$$h_t(x) = \frac{1}{\sqrt{t!}} \sum_{a=0}^{t/2} (-1)^a x^{t-2a} \frac{t!}{(t-2a)!a!2^a} \,.$$

If $|x| = O(\sqrt{t})$, the largest term above will occur when $x^2a \approx (t - 2a)^2$, which occurs when a is a constant multiple of t. In this case, it is easy to see that the term in question is $O(1)^t$. Thus, when x and x' come from the same component, $P_t = O(s^2/t)^t$ (with a small enough implied constant in the big-O) with high probability.

Therefore, by estimating P_t to error $O(s)^t$, we can reliably distinguish between the cases where x and x' come from the same component and the one where they do not. This can be done directly using Proposition 1.7 along with Equation (9) and the fact that $H_t((x - x')/\sqrt{2}) =$ $\mathbf{E}_{y \sim N(0,I)}[H_n((x - x')/\sqrt{2}, y)]$. Applying Lemma 4.3, we find that both of these estimators have second moment bounded by $O(1 + s^2/t)^t$. This can be estimated to appropriate error with $poly(k/(w_{\min}\tau))$ samples.

4.4 Learning Mixtures of Linear Regressions

Problem Setup We start by defining the underlying probabilistic model. A linear regression problem produces a distribution on \mathbb{R}^{d+1} . In particular, given $\sigma > 0$ and $\beta \in \mathbb{R}^d$, we define a distribution $L_{\beta,\sigma}$ on \mathbb{R}^{d+1} given as the distribution of (X, y), where $X \in \mathbb{R}^d$ is distributed as N(0, I) and $y = \beta \cdot X + N(0, \sigma^2)$. Notice that as a distribution over \mathbb{R}^{d+1} , $L_{\beta,\sigma}$ is just a Gaussian, namely

$$L_{\beta,\sigma} \sim N\left(0, \begin{bmatrix} I & \beta \\ \beta^T & \|\beta\|_2^2 + \sigma^2 \end{bmatrix}\right).$$

We can now define the corresponding mixture model.

Definition 4.17 (Mixtures of Linear Regressions). A k-mixture of linear regressions (k-MLR) with error $\sigma > 0$, is any distribution of the form $F = \sum_{i=1}^{k} w_i L_{\beta_i,\sigma}$ for some unknown mixing weights $w_i \ge 0$ with $\sum_{i=1}^{k} w_i = 1$ and unknown vectors β_i . We will assume that the β_i 's have some known upper bound B on their norms, namely $\|\beta_i\|_2 \le B$ for $1 \le i \le k$.

Definition 4.18 (Density Estimation for Mixtures of Linear Regressions). The density estimation problem for mixtures of linear regressions is the following: The input is a multiset of i.i.d. samples in \mathbb{R}^{d+1} drawn from an unknown k-mixture $F = \sum_{i=1}^{k} w_i L_{\beta_i,\sigma}$, where σ, k , and B are known. The goal of the learner is to output a (sampler for a) hypothesis distribution H such that with high probability $d_{\text{TV}}(H, F) \leq \epsilon$.

Our main algorithmic result for this problem is the following:

Theorem 4.19 (Density Estimation Algorithm for k-MLR). Suppose that we are given sample access to a k-MLR distribution F with $B, \sigma \leq 1$. Then there exists an algorithm that given k, σ , and $\epsilon > 0$, draws $N = \text{poly}(k, d)\epsilon^{O(\sigma^{-2})}$ samples from F, runs in poly(N, d) time, and returns a sampler for a distribution that is ϵ -close to F in total variation distance.

Note that if we are given an F with B or σ more than 1, we can simply divide the *y*-values by some constant C, which will have the effect of dividing both B and σ by C. For C large enough, the hypotheses of Theorem 4.19 should apply.

Proof of Theorem 4.19 By dividing the *y*-values by 2, we get a new distribution of the same form with σ and β_i all half as big. Thus, we can reduce to the case where $\sigma, B \leq 1/2$, which we assume below.

As in the Gaussian mixtures application, we will produce this distribution through rejection sampling. In particular, we will produce a random sample x from the standard Gaussian G, and

then we will try to accept it with probability proportional to (F/G)(x). This depends on being able to approximate the function (F/G)(x) with small L_1 error. We do this by trying to compute its Taylor expansion using Proposition 1.7.

In the lemma below, we start by defining the appropriate parameter moments for this setting.

Lemma 4.20. Writing F = (X, y) with $X \in \mathbb{R}^d$ and $y \in \mathbb{R}$, we have that

$$M_m := \sum_{i=1}^k w_i \beta_i^{\otimes m} = \mathbf{E}_{(X,y) \sim F, Y \sim N(0,I)} \left[(y^m / \sqrt{m!}) H_m(X,Y) \right].$$
(10)

Furthermore, if $B, \sigma \leq 1$, the second moment of $(y^m/\sqrt{m!})H_n(X,Y)$ above in any direction is bounded by

$$V = 2^{O(m)}.$$

Proof. To begin with, we note that the expectation of Y of the right hand side of Equation (10) is $\mathbf{E}_{(X,y)\sim F}[(y^m/\sqrt{m!})H_m(X)]$. Since F is a mixture of the $L_{\beta_i,\sigma}$, this is

$$\sum_{i=1}^k w_i \mathbf{E}_{(X,y)\sim L_{\beta_i,\sigma}}[(y^m/\sqrt{m!})H_m(X)].$$

For the i^{th} term of this sum, we note that $y = \beta_i \cdot X + \xi$, where $\xi \sim N(0, \sigma^2)$ is an independent Gaussian. Thus, $y^m / \sqrt{m!} = (\beta_i \cdot X + \xi)^m / \sqrt{m!}$. For any fixed value of ξ , this means that $(y^m / \sqrt{m!})$ is $\langle \beta^{\otimes m}, H_m(X) \rangle$ plus a polynomial of degree less than m in X. As $H_m(X)$ is orthogonal to polynomials of degree less than m, the above expectation is

$$\sum_{i=1}^{k} w_i \mathbf{E}_X[\langle \beta_i^{\otimes m}, H_m(X) \rangle, H_m(x)] = M_m \,,$$

by the fact that $\beta_i^{\otimes m}$ is symmetric and the orthonormality properties of the Hermite tensors H_m .

For the second moment bound, we let T be a tensor with $||T||_2 \leq 1$ and we wish to bound

$$\mathbf{E}_{(X,y)\sim F,Y\sim N(0,I)}[((y^m/\sqrt{m!})\langle T,H_m(X,Y)\rangle)^2].$$

We can do this by bounding the expectation on each component of the mixture. We note that with $(X, y) \sim L_{\beta_i}$ and $Y \sim N(0, I)$, all of the terms in the above expectation are polynomials of a Gaussian input. By Lemma 4.3, we have that $\|\langle T, H_m(X, Y) \rangle\|_2 = O(1)^m$. Similarly, a direct computation shows that $\|(y^m/\sqrt{m!})\|_2 = O(1)^m$. Applying Gaussian Hypercontractivity (Fact 2.4) and Holder's Inequality, we conclude that $\|(y^m/\sqrt{m!})\langle T, H_m(X, Y) \rangle\|_2 = O(1)^m$, as desired. \Box

We next need to determine how to approximate (F/G)(x). Using Hermite analysis, we know that it is given by

$$(F/G)(x) = \sum_{n=0}^{\infty} \langle T_n, H_n(x) \rangle ,$$

where

$$T_n = \mathbf{E}[H_n(F)] = \sum_{i=1}^k w_i \mathbf{E}[H_n(L_{\beta_i,\sigma})].$$

As $L_{\beta,\sigma}$ is just a Gaussian, we have by Lemma 2.7 of [Kan21] (noting the difference in normalization between their Hermite polynomials and ours), we have that $\mathbf{E}[H_n(L_{\beta,\sigma})]$ is 0 if n is odd and is

$$\frac{(n-1)!!}{\sqrt{n!}} \operatorname{Sym}\left(\begin{bmatrix} 0 & \beta \\ \beta^T & \|\beta\|_2^2 + \sigma^2 - 1 \end{bmatrix}^{\otimes n/2} \right)$$

if n is even. In particular, letting e_y be the unit vector in the y-direction, this is

$$\frac{(n-1)!!}{\sqrt{n!}}\operatorname{Sym}\left(\left(2\beta \otimes e_y + \|\beta\|_2^2 e_y \otimes e_y + (\sigma^2 - 1)e_y \otimes e_y\right)^{\otimes n/2}\right)$$
$$= \frac{(n-1)!!}{\sqrt{n!}} \sum_{a+b+c=n/2} \binom{n/2}{a,b,c} 2^a \|\beta\|_2^{2b} (\sigma^2 - 1)^c \operatorname{Sym}\left(\beta^{\otimes a} \otimes e_y^{\otimes n-a}\right)$$

Therefore, we have that

$$\begin{split} \langle T_n, H_n(X) \rangle &= \sum_{i=1}^{\kappa} w_i \langle \mathbf{E}[H_n(L_{\beta_i,\sigma})], H_n(X) \rangle \\ &= \frac{(n-1)!!}{\sqrt{n!}} \sum_{i=1}^{k} w_i \sum_{a+b+c=n/2} \binom{n/2}{a, b, c} 2^a ||\beta||_2^{2b} (\sigma^2 - 1)^c \langle \beta^{\otimes a} \otimes e_y^{\otimes n-a}, H_n(X) \rangle \\ &= \frac{(n-1)!!}{\sqrt{n!}} \sum_{a+b+c=n/2} \binom{n/2}{a, b, c} 2^a (\sigma^2 - 1)^c \langle M_{a+2b} \otimes e_y^{\otimes n-a}, H_n(X) \otimes I_d^{\otimes b} \rangle \\ &= \frac{(n-1)!!}{\sqrt{n!}} \sum_{a+b+c=n/2} \binom{n/2}{a, b, c} 2^a (\sigma^2 - 1)^c \mathbf{E}_Y \left[\langle M_{a+2b} \otimes e_y^{\otimes n-a}, H_n(X, Y) \otimes I_d^{\otimes b} \rangle \right] \,. \end{split}$$

Using Proposition 1.7 and Lemmas 4.20 and 4.3, we can approximate $\langle M_{a+2b} \otimes e_y^{\otimes n-a}, H_n(X,Y) \otimes I_d^{\otimes b} \rangle$ to expected squared error $O(1)^n / \sqrt{N}$ with N samples in poly(N, k, d) time. As the above sum has poly(n) terms with factors $O(1)^n$ (so long as $\sigma \leq 1$), we can approximate $\langle T_n, H_n(X) \rangle$ for random Gaussian X to expected error $O(1)^n / \sqrt{N}$ in polynomial time.

Unfortunately, we cannot compute infinitely many values of $\langle T_n, H_n(x) \rangle$, so we will want to show that we can truncate the sum. This amounts to showing that $||T_n||_2$ is small when n is sufficiently large. For this, we note that

$$T_n = \frac{(2n-1)!!}{\sqrt{n!}} \sum_{i=1}^k w_i \text{Sym}\left((2\beta_i \otimes e_y + (\|\beta_i\|_2^2 + \sigma^2 - 1)e_y \otimes e_y)^{\otimes n/2} \right)$$

First, using the triangle inequality, we note that

$$||T_n||_2 \le \max_i \left\| \operatorname{Sym}\left((2\beta_i \otimes e_y + (||\beta_i||_2^2 + \sigma^2 - 1)e_y \otimes e_y)^{\otimes n/2} \right) \right\|_2.$$

Next we note that the matrix $(2\beta_i \otimes e_y + (||\beta_i||_2^2 + \sigma^2 - 1)e_y \otimes e_y)$ is similar to some 2×2 matrix

 $M_i = \begin{bmatrix} \gamma_{1,i} & 0\\ 0 & \gamma_{2,i} \end{bmatrix}$. Thus, we have that

$$\begin{split} \|T_n\|_2^2 &\leq \max_i \left\| \operatorname{Sym}(M_i^{\otimes n/2}) \right\|_2^2 \\ &= \max_i \left\| \sum_{a=0}^{n/2} \binom{n/2}{a} \gamma_{1,i}^a \gamma_{2,i}^{n/2-a} \operatorname{Sym}(e_1^{\otimes 2a} \otimes e_2^{\otimes n-2a}) \right\|_2^2 \\ &= \max_i \sum_{a=0}^{n/2} \binom{n/2}{a}^2 \gamma_{1,i}^{2a} \gamma_{2,i}^{n-2a} \binom{n}{2a}^{-1}, \end{split}$$

where the last line is because symmetrizations of $e_1^{\otimes 2a} \otimes e_2^{\otimes n-2a}$ for different values of a are orthogonal, and two permutations of $e_1^{\otimes 2a} \otimes e_2^{\otimes n-2a}$ are orthogonal unless the copies of e_1 end up in the same places (in which case they have dot product 1). Thus, we conclude that

$$||T_n||_2^2 \le n \max_{i,a} \gamma_{1,i}^{2a} \gamma_{2,i}^{n-2a} \le n \gamma^n$$
,

where γ is the largest absolute value of any eigenvalue of any of the matrices $(2\beta_i \otimes e_y + (||\beta_i||_2^2 + \sigma^2 - 1)e_y \otimes e_y)$.

To bound the latter quantity, note that M_i has determinant $-\|\beta_i\|_2^2$ and trace $\|\beta_i\|_2^2 + \sigma^2 - 1$. Therefore, it has eigenvalues

$$\frac{\|\beta_i\|_2^2 + \sigma^2 - 1 \pm \sqrt{(\|\beta_i\|_2^2 + \sigma^2 - 1)^2 + 4\|\beta_i\|_2^2}}{2}$$

Since $(\|\beta_i\|_2^2 + \sigma^2 - 1)^2 + 4\|\beta_i\|_2^2 \le (1 + \|\beta_i\|_2^2)^2$, these eigenvalues are between

$$\frac{\|\beta_i\|_2^2 + \sigma^2 - 1 \pm (1 + \|\beta_i\|_2^2)}{2}$$

or between $-1 + \sigma^2/2$ and $\|\beta_i\|_2^2 + \sigma^2/2$. In particular, if $B, \sigma \leq 1/2$ we have that $\gamma \leq 1 - \sigma^2/2$.

Therefore, if this holds, up to L_2 error $\epsilon/2$, we have that

$$(F/G)(x) = \sum_{n=0}^{O(\sigma^{-2}\log(1/\epsilon))} \langle T_n, H_n(x) \rangle$$

which can be approximated to error $\epsilon/2$ using $N = \text{poly}(k, d)\epsilon^{O(\sigma^{-2})}$ samples and poly(N, d) time.

Finally, noting that F is a mixture of Gaussians with covariances bounded below by $\Omega(\sigma^2)$ and above by a constant, it is not hard to see that except with ϵ probability, a point sampled from F has F/G(x) at most poly $(1/(\epsilon\sigma))$. Thus, if C is a big enough polynomial in $1/(\epsilon\sigma)$, then picking a random sample $x \sim G$ and accepting with probability min(1, f(x)/C), where f(x) is our approximation to (F/G)(x), gives an ϵ -approximation to the distribution F.

This concludes the proof of Theorem 4.19.

5 Conclusions

In this work, we gave a general efficient algorithm to approximate higher-order moment tensors, as long as there are reasonable unbiased estimators for these moments that can be efficiently sampled. This type of tensors arise in a range of learning problems. We leveraged our general result to obtain the first polynomial-time algorithms for learning mixtures of linear regressions, learning sums of ReLUs, density estimation for mixtures of spherical Gaussians with bounded means, and parameter estimation for mixtures of spherical Gaussians under optimal separation. In all cases, our learning algorithms run in poly(d, k) time.

A number of open questions remain. The most obvious ones are quantitative: Specifically, can the dependence on $1/\epsilon$ for learning sums of ReLUs be improved to polynomial? Can we remove the bounded means assumption for density estimation of mixtures of spherical Gaussians? Can we obtain a polynomial-time parameter estimation algorithm for mixtures of linear regressions? More broadly, for what other learning tasks is the implicit moment tensor technique applicable? These are left as interesting directions for future work.

References

- [ABH⁺18] H. Ashtiani, S. Ben-David, N. J. A. Harvey, C. Liaw, A. Mehrabian, and Y. Plan. Nearly tight sample complexity bounds for learning mixtures of gaussians via sample compression schemes. In Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, pages 3416–3425, 2018.
- [AK01] S. Arora and R. Kannan. Learning mixtures of arbitrary Gaussians. In Proceedings of the 33rd Symposium on Theory of Computing, pages 247–257, 2001.
- [AM05] D. Achlioptas and F. McSherry. On spectral learning of mixtures of distributions. In Proceedings of the Eighteenth Annual Conference on Learning Theory (COLT), pages 458–469, 2005.
- [BJW19] A. Bakshi, R. Jayaram, and D. P. Woodruff. Learning two layer rectified neural networks in polynomial time. In *Conference on Learning Theory*, COLT 2019, pages 195–268, 2019.
- [Bon70] A. Bonami. Etude des coefficients fourier des fonctiones de $l^p(g)$. Ann. Inst. Fourier (Grenoble), 20(2):335–402, 1970.
- [BRST21] J. Bruna, O. Regev, M. J. Song, and Y. Tang. Continuous LWE. In STOC '21: 53rd Annual ACM SIGACT Symposium on Theory of Computing, 2021, pages 694–707. ACM, 2021.
- [BV08] S. C. Brubaker and S. Vempala. Isotropic PCA and Affine-Invariant Clustering. In *Proc.* 49th IEEE Symposium on Foundations of Computer Science, pages 551–560, 2008.
- [CDG⁺23] S. Chen, Z. Dou, S. Goel, A. R. Klivans, and R. Meka. Learning narrow one-hiddenlayer relu networks. CoRR, abs/2304.10524, 2023. Conference version in COLT'23.
- [CGKM22] S. Chen, A. Gollakota, A. R. Klivans, and R. Meka. Hardness of noise-free learning for two-hidden-layer neural networks. In *NeurIPS*, 2022.

- [Che20] S. Chen. Learning mixtures of linear regressions in subexponential time via fourier moments. Conference talk at STOC 2020, 2020. Available at https://www.youtube.com/watch?v=ed_Vz-cDlz4.
- [CKM21] S. Chen, A. R. Klivans, and R. Meka. Learning deep relu networks is fixed-parameter tractable. In 62nd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2021, pages 696–707. IEEE, 2021.
- [CLS20] S. Chen, J. Li, and Z. Song. Learning mixtures of linear regressions in subexponential time via fourier moments. In *Proceedings of the 52nd Annual ACM SIGACT Symposium* on Theory of Computing, STOC 2020, pages 587–600. ACM, 2020.
- [CN24] S. Chen and S. Narayanan. A faster and simpler algorithm for learning shallow networks. In *The Thirty Seventh Annual Conference on Learning Theory*, volume 247 of *Proceedings of Machine Learning Research*, pages 981–994. PMLR, 2024.
- [Das99] S. Dasgupta. Learning mixtures of Gaussians. In Proceedings of the 40th Annual Symposium on Foundations of Computer Science, pages 634–644, 1999.
- [DeV89] R. D. DeVeaux. Mixtures of linear regressions. Computational Statistics & Data Analysis, 8(3):227–245, November 1989.
- [DFS16] A. Daniely, R. Frostig, and Y. Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. In Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, pages 2253–2261, 2016.
- [DK20] I. Diakonikolas and D. M. Kane. Small covers for near-zero sets of polynomials and learning latent variable models. In 61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, pages 184–195, 2020. Full version available at https://arxiv.org/abs/2012.07774.
- [DK24] I. Diakonikolas and D. M. Kane. Efficiently learning one-hidden-layer relu networks via Schur polynomials. In The Thirty Seventh Annual Conference on Learning Theory, volume 247 of Proceedings of Machine Learning Research, pages 1364–1378. PMLR, 2024.
- [DKK⁺16] I. Diakonikolas, G. Kamath, D. M. Kane, J. Li, A. Moitra, and A. Stewart. Robust estimators in high dimensions without the computational intractability. In *Proceedings* of FOCS'16, pages 655–664, 2016. Journal version in SIAM Journal on Computing, 48(2), pages 742-864, 2019.
- [DKKZ20] I. Diakonikolas, D. M. Kane, V. Kontonis, and N. Zarifis. Algorithms and SQ lower bounds for PAC learning one-hidden-layer relu networks. In *Conference on Learning Theory, COLT 2020*, volume 125 of *Proceedings of Machine Learning Research*, pages 1514–1539. PMLR, 2020.
- [DKPP24] I. Diakonikolas, S. Karmalkar, S. Pang, and A. Potechin. Sum-of-squares lower bounds for non-gaussian component analysis. In 65th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2024, 2024.

- [DKPZ23] I. Diakonikolas, D. M. Kane, T. Pittas, and N. Zarifis. SQ lower bounds for learning mixtures of separated and bounded covariance gaussians. In *The Thirty Sixth Annual Conference on Learning Theory, COLT 2023*, volume 195 of *Proceedings of Machine Learning Research*, pages 2319–2349. PMLR, 2023.
- [DKS17] I. Diakonikolas, D. M. Kane, and A. Stewart. Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. In 58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, pages 73–84, 2017. Full version at http://arxiv.org/abs/1611.03473.
- [DKS18] I. Diakonikolas, D. M. Kane, and A. Stewart. List-decodable robust mean estimation and learning mixtures of spherical gaussians. In *Proceedings of the 50th Annual ACM* SIGACT Symposium on Theory of Computing, STOC 2018, pages 1047–1060, 2018. Full version available at https://arxiv.org/abs/1711.07211.
- [FOS06] J. Feldman, R. O'Donnell, and R. Servedio. PAC learning mixtures of Gaussians with no separation assumption. In *COLT*, pages 20–34, 2006.
- [GGJ⁺20] S. Goel, A. Gollakota, Z. Jin, S. Karmalkar, and A. R. Klivans. Superpolynomial lower bounds for learning one-layer neural networks using gradient descent. In *Proceedings* of the 37th International Conference on Machine Learning, ICML 2020, volume 119 of Proceedings of Machine Learning Research, pages 3587–3596, 2020.
- [GK19] S. Goel and A. R. Klivans. Learning neural networks with two nonlinear layers in polynomial time. In *Conference on Learning Theory*, COLT 2019, pages 1470–1499, 2019.
- [GKKT17] S. Goel, V. Kanade, A. R. Klivans, and J. Thaler. Reliably learning the relu in polynomial time. In *Proceedings of the 30th Conference on Learning Theory*, COLT 2017, pages 1004–1042, 2017.
- [GKLW19] R. Ge, R. Kuditipudi, Z. Li, and X. Wang. Learning two-layer neural networks with symmetric inputs. In 7th International Conference on Learning Representations, ICLR 2019, 2019.
- [GLM18] R. Ge, J. D. Lee, and T. Ma. Learning one-hidden-layer neural networks with landscape design. In 6th International Conference on Learning Representations, ICLR 2018, 2018.
- [Gro75] L. Gross. Logarithmic Sobolev inequalities. Amer. J. Math., 97(4):1061–1083, 1975.
- [GVV22] A. Gupte, N. Vafa, and V. Vaikuntanathan. Continuous LWE is as hard as LWE & applications to learning gaussian mixtures. In 63rd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2022, 2022, pages 1162–1173. IEEE, 2022.
- [HL18] S. B. Hopkins and J. Li. Mixture models, robustness, and sum of squares proofs. In Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, pages 1021–1034, 2018.
- [HP15] M. Hardt and E. Price. Tight bounds for learning a mixture of two gaussians. In Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, pages 753–760, 2015.

- [JJ94] M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the EM algorithm. Neural Computation, 6(2):181–214, 1994.
- [JSA15] M. Janzamin, H. Sedghi, and A. Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods, 2015.
- [Kan21] D. M. Kane. Robust learning of mixtures of gaussians. In Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, 2021, pages 1246–1258. SIAM, 2021. Also available at https://arxiv.org/abs/2007.05912.
- [KC19] J. Kwon and C. Caramanis. EM converges for a mixture of many linear regressions. CoRR, abs/1905.12106, 2019.
- [KSS18] P. K. Kothari, J. Steinhardt, and D. Steurer. Robust moment estimation and improved clustering via sum of squares. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018*, pages 1035–1046, 2018.
- [KSV08] R. Kannan, H. Salmasian, and S. Vempala. The spectral method for general mixture models. SIAM J. Comput., 38(3):1141–1156, 2008.
- [Liu22] A. Liu. Clustering mixtures with almost optimal separation in polynomial time. Conference talk at STOC 2022, 2022. Available at https://www.youtube.com/watch?v=xkU5gSHGar4.
- [Liu24] A. Liu. Clustering mixtures with almost optimal separation in polynomial time. Invited talk at Workshop on New Frontiers in Robust Statistics, TTI-Chicago, June 2024, 2024. Available at http://www.iliasdiakonikolas.org/ttic-robust24/Liu.pdf.
- [LL18] Y. Li and Y. Liang. Learning mixtures of linear regressions with nearly optimal complexity. In Conference On Learning Theory, COLT 2018, volume 75 of Proceedings of Machine Learning Research, pages 1125–1144. PMLR, 2018.
- [LL22] A. Liu and J. Li. Clustering mixtures with almost optimal separation in polynomial time. In STOC '22: 54th Annual ACM SIGACT Symposium on Theory of Computing, 2022, pages 1248–1261, 2022. Full version available at https://arxiv.org/abs/2112.00706.
- [LS17] J. Li and L. Schmidt. Robust and proper learning for mixtures of gaussians via systems of polynomial inequalities. In *Proceedings of the 30th Conference on Learning Theory, COLT 2017*, volume 65 of *Proceedings of Machine Learning Research*, pages 1302–1382. PMLR, 2017.
- [LZZ24] S. Li, I. Zadik, and M. Zampetakis. On the hardness of learning one hidden layer neural networks. *CoRR*, abs/2410.03477, 2024. Available at https://arxiv.org/abs/2410.03477.
- [MV10] A. Moitra and G. Valiant. Settling the polynomial learnability of mixtures of Gaussians. In *FOCS*, pages 93–102, 2010.
- [Pea94] K. Pearson. Contribution to the mathematical theory of evolution. Phil. Trans. Roy. Soc. A, 185:71–110, 1894.
- [RV17] O. Regev and A. Vijayaraghavan. On learning mixtures of well-separated gaussians. In 58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, pages 85–96, 2017.

- [SJA16] H. Sedghi, M. Janzamin, and A. Anandkumar. Provable tensor methods for learning mixtures of generalized linear models. In *Proceedings of the 19th International Confer*ence on Artificial Intelligence and Statistics, AISTATS 2016, pages 1223–1231, 2016.
- [SOAJ14] A. T. Suresh, A. Orlitsky, J. Acharya, and A. Jafarpour. Near-optimal-sample estimators for spherical Gaussian mixtures. In NIPS 2014, pages 1395–1403, 2014.
- [VW02] S. Vempala and G. Wang. A spectral algorithm for learning mixtures of distributions. In Proceedings of the 43rd Annual Symposium on Foundations of Computer Science, pages 113–122, 2002.
- [VW19] S. Vempala and J. Wilmes. Gradient descent for one-hidden-layer neural networks: Polynomial convergence and SQ lower bounds. In *Conference on Learning Theory*, *COLT 2019*, pages 3115–3117, 2019.
- [ZJD16] K. Zhong, P. Jain, and I. S. Dhillon. Mixed linear regression with multiple components. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, pages 2190–2198, 2016.
- [ZLJ16] Y. Zhang, J. D. Lee, and M. I. Jordan. L1-regularized neural networks are improperly learnable in polynomial time. In *Proceedings of the 33nd International Conference on Machine Learning*, *ICML 2016*, pages 993–1001, 2016.
- [ZSJ⁺17] K. Zhong, Z. Song, P. Jain, P. L. Bartlett, and I. S. Dhillon. Recovery guarantees for one-hidden-layer neural networks. In *Proceedings of the 34th International Conference* on Machine Learning, ICML 2017, pages 4140–4149, 2017.