# Learning Mixtures of Gaussians Using Diffusion Models

Khashayar Gatmiry
gatmiry@mit.edu
MIT

Jonathan Kelner
kelner@mit.edu
MIT

Holden Lee
hlee283@jhu.edu
JHU

April 30, 2024

## Abstract

We give a new algorithm for learning mixtures of $k$ Gaussians (with identity covariance in $\mathbb{R}^n$) to TV error $\varepsilon$, with quasi-polynomial ($O(n^{\operatorname{poly\,log}\left(\frac{n+k}{\varepsilon}\right)})$)) time and sample complexity, under a minimum weight assumption. Unlike previous approaches, most of which are algebraic in nature, our approach is analytic and relies on the framework of diffusion models. Diffusion models are a modern paradigm for generative modeling, which typically rely on learning the score function (gradient log-pdf) along a process transforming a pure noise distribution, in our case a Gaussian, to the data distribution. Despite their dazzling performance in tasks such as image generation, there are few end-to-end theoretical guarantees that they can efficiently learn nontrivial families of distributions; we give some of the first such guarantees. We proceed by deriving higher-order Gaussian noise sensitivity bounds for the score functions for a Gaussian mixture to show that that they can be inductively learned using piecewise polynomial regression (up to poly-logarithmic degree), and combine this with known convergence results for diffusion models. Our results extend to continuous mixtures of Gaussians where the mixing distribution is supported on a union of $k$ balls of constant radius. In particular, this applies to the case of Gaussian convolutions of distributions on low-dimensional manifolds, or more generally sets with small covering number.

# Contents

# 1  Introduction

We address the problem of learning a mixture of $k$ Gaussians (with identity covariance) from samples, using the framework of diffusion models. Our motivation is twofold: First, as one of the simplest but nevertheless challenging mixture models, this is a classic learning problem in statistics and computer science. Second, diffusion models are a empirically successful paradigm for generative modeling, which work well for learning multimodal distributions in practice but for which theoretical guarantees are lacking. By applying diffusion models to the problem of learning Gaussian mixtures, we obtain a new method distinct from the common algebraic approaches, as well as give theoretical grounding to the success of diffusion models.

Formally, we wish to learn the following distribution on $\mathbb{R}^n$ from iid samples:

$$P_0 = Q_0 * \mathcal{N}(0, \sigma_0^2 I_n),$$

which we think of as a (possibly continuous) mixture of Gaussians, where $Q_0$ is the distribution of the means. In the special case that $Q_0 = \sum_{j=1}^k \alpha_j \delta_{\mu_j}$, this is exactly a mixture of $k$ Gaussians; however, our results hold in this more general setting. Note that if the covariance is known but non-identity, we can first transform the data to be in this setting. Our goal is distribution learning: to output samples from a distribution $\varepsilon$-close in TV distance to the actual one.

For $x_0 \in \mathbb{R}^n$, let $B_R(x_0) = \{x \in \mathbb{R}^n : |x - x_0| \leq R\}$ denote the closed ball of radius $R$ around $x$. We make the following assumptions.

**Assumption 1.1:** Fix $R_0 \geq 1$, $D$, and $k$. The following hold:

1. For every point $\mu$ in the support of $Q_0$, we have $Q_0(B_{R_0}(\mu)) \geq \alpha_{\min}$.

2. There exist $\overline{\mu}_1, \ldots, \overline{\mu}_k$ such that the support of $Q_0$ is contained in $\bigcup_{i=1}^k B_{R_0}(\overline{\mu}_i)$.

3. $Q_0(B_D(0)) = 1$.

This is a strict generalization of a mixture of $k$ Gaussians: we allow a mixture of $k$ arbitrary distributions supported on balls of radius $R_0$ convolved with Gaussians. Our main theorem is that these mixtures can be learned with quasi-polynomial time and samples with an algorithm based on diffusion models.

**Theorem 1.2.** *Given $\varepsilon > 0$ with $\varepsilon \leq \min\left\{\frac{1}{2}, \frac{\sigma_0}{R_0}, \frac{1}{D}, \frac{1}{n}, \alpha_{\min}\right\}$, and given Assumption 1.1, Algorithm 1 learns a distribution that is $\varepsilon$-close in TV distance to $P_0$ with time and sample complexity $\left(n \ln\left(\frac{1}{\delta}\right)\right)^{O\left(\left(\ln\left(\frac{1}{\varepsilon}\right)^3 + \left(\frac{R_0}{\sigma_0}\right)^6\right) \ln\left(\frac{1}{\varepsilon}\right)^4\right)}$ with probability $\geq 1 - \delta$.*

Note that because of the restriction on $\varepsilon$, $\ln\left(\frac{1}{\varepsilon}\right)$ implicitly has logarithmic dependence on $D$, $n$, and $\frac{1}{\alpha_{\min}} \geq k$. In the case where $\frac{R_0}{\sigma_0}$ is a constant, we can remove that dependence. We further remark that in the case of a (discrete) mixture of $k$ Gaussians, a straightforward SVD pre-processing step can replace the dependence on $n$ to $\min\{n, k\}$, at an extra additive cost polynomial in $n$ (see e.g., [VW04]).

Algorithm 1 consists of two parts: The "learning" part involves learning the score functions of distributions that bridge the data distribution with a pure noise (Gaussian) distribution. Once these scores are obtained, the "generation" part uses these learned score functions and can generate

---

**Algorithm 1** Learning Gaussian mixture with diffusion model

---
1: **Input:** Error $\varepsilon$, failure probability $\delta$, sample access to mixture $Q_0$ satisfying Assumption 1.1.
     $\triangleright$ Learning
2: Let $t_1 = \frac{\varepsilon^2 \sigma_0^2}{2\sqrt{n}}$ and choose a step size schedule $t_1 < \cdots < t_N$ as in Theorem 3.1.
3: Let $d = \Theta\left(\left(\ln\left(\frac{1}{\varepsilon}\right)^3 + \left(\frac{R_0}{\sigma_0}\right)^6\right)\ln\left(\frac{1}{\varepsilon}\right)^4\right)$ for an appropriate constant.
4: Set $\mathcal{C}_{N_{\text{step}}} = \{0\}$.                                              $\triangleright$ Set of warm starts
5: **for** $\ell$ from $N_{\text{step}}$ to 1 **do**
6:     Let $V_1, \ldots, V_{k_\ell}$ be the Voronoi partition induced by $\mathcal{C}_\ell = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_{k_\ell}\}$.
7:     Draw $N = \left(n \ln\left(\frac{1}{\delta}\right)\right)^{\Theta(d)}$ samples $x_1, \ldots, x_N \sim P_0$ and let $y_i = x_i + \sqrt{t_\ell} \cdot \xi_i$, $\xi_i \sim \mathcal{N}(0, I_n)$.
8:     **for** $j$ from 1 to $k_\ell$ **do**
9:        Let $\sigma_\ell^2 = t_\ell + \sigma_0^2$ and

$$(\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \le d} = \underset{(b_{\mathbf{k}})_{|\mathbf{k}| \le d} \in B_D(0)}{\arg\min} \sum_{i: y_i \in V_j} \left| \sum_{|\mathbf{k}| \le d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i - \widehat{\mu}_j) - \left(\left(1 - \frac{\sigma_\ell^2}{t_\ell}\right) y_i + \frac{\sigma_\ell^2}{t_\ell} x_i\right) \right|^2$$

10:     **end for**
11:     Define $\widehat{g}_\ell(y) := \sum_{i=1}^{k_\ell} \mathbb{1}_{V_i}(y) \sum_{|\mathbf{k}| \le d} \widehat{b}_{\mathbf{k}}^{(i)} h_{\mathbf{k}}(y - \widehat{\mu}_i)$ and $s_\ell(y) = \frac{\widehat{g}_\ell(y) - y}{\sigma_\ell^2}$.
12:     **if** time has halved since last computation of warm starts **then**
13:        Run Algorithm 2 with fresh samples, and failure probability $\frac{\delta}{2N_{\text{step}}}$ to obtain $\mathcal{C}_{\ell-1}$.
14:     **else** let $\mathcal{C}_{\ell-1} = \mathcal{C}_\ell$
15:     **end if**
16: **end for**
17: **Output:** Score functions $s_1, \ldots, s_{N_{\text{step}}}$.
18: **Input:** Score functions $s_1, \ldots, s_{N_{\text{step}}}$ for times $t_1 < \cdots < t_{N_{\text{step}}}$          $\triangleright$ Generation
19: Draw $\widehat{y}_{t_{N_{\text{step}}}} \sim \mathcal{N}(0, t_{N_{\text{step}}} I_n)$.
20: **for** $\ell$ from $N$ to 2 **do**
21:     Let $\widehat{y}_{t_{\ell-1}} = \widehat{y}_{t_\ell} + 2\left[(t_\ell - 1) - \sqrt{(t_{\ell-1} - 1)(t_\ell - 1)}\right] s_{t_\ell}(\widehat{y}_{t_\ell}) + \sqrt{h_{\ell-1}} \cdot \xi_{t_\ell}$ where $\xi_{t_\ell} \sim \mathcal{N}(0, I_n)$.
22: **end for**
23: **Output:** $\widehat{y}_{t_1}$

---

as many samples as desired. The high-probability bound is over the learning part: with high probability, the learned score functions are such that the generation process satisfies the TV distance bound.

A special case of Theorem 1.2 is the problem of learning a distribution that is equal to a distribution on a low-dimensional manifold convolved with a Gaussian. The "manifold" assumption that we need is much weaker: simply that it has can be covered by $C^d$ balls of radius $R_0$, for some constant $C$. It is straightforward to obtain the following.

**Corollary 1.3.** *Fix $\sigma_0 = 1$ and constants $R_0, C > 1$. Let $0 < \varepsilon < \frac{1}{2}$. Suppose that $Q_0$ is supported on a set $M$ such that $M$ has radius $D$, $M$ can be covered with $C^d$ balls of radius $R_0$, and such that for every point $\mu$ in the support of $Q_0$, $Q_0(B_{R_0}(\mu)) \ge \frac{\varepsilon}{C^d}$. Then Algorithm 1 learns a distribution that is $\varepsilon$-close in TV distance to $P_0$ with time and sample complexity $\left(n \ln\left(\frac{1}{\delta}\right)\right)^{\left(d + \ln\left(\frac{nD}{\varepsilon}\right)\right)^7}$ with*

4

*probability* $\geq 1 - \delta$.

We note that this is a setting where diffusion models can provably learn under a manifold assumption, but in contrast to most prior work, the distribution *cannot* be learned using straightforward methods such as through binning or kernel density estimation. For example, if $M$ can be covered with $\left(\frac{C}{\varepsilon}\right)^d$ balls of radius $\varepsilon$, then learning a distribution exactly supported on $M$ can be done to Wasserstein distance $\varepsilon$ with complexity $\widetilde{O}\left(\left(\frac{C}{\varepsilon}\right)^d\right)$ simply with a binning procedure. However, we consider learning a distribution on $M$ convolved with a Gaussian, which is a more challenging problem.

We note that while Theorem 1.2 is known in the case of a mixture of $k$ Gaussians using a completely different algorithm based on algebraic methods [DK20] (without the requirement on $\alpha_{\min}$ and with better exponents), to our knowledge the extension to our "generalized" mixture of $k$ Gaussians is new. It is interesting to note that the generic framework of diffusion models allows us to match (up to polylogarithmic factors in the exponent) guarantees obtainable using more specialized algebraic procedures.

Note that we do *not* proceed by learning the density function; instead, "learning" the distribution means that we have a procedure to generate an additional sample, each step of which involves evaluation of a learned score function. It may be possible to upgrade this guarantee to a guarantee of learning the density [QR23].

We leave open the questions of removing the requirement on $\alpha_{\min}$ and improving the polynomial dependence on $\ln\left(\frac{1}{\varepsilon}\right)$ (in Theorem 1.2) and $d$ (in Theorem 1.3). The question remains of whether the complexity of learning mixtures of Gaussians (of equal known covariance) is truly quasi-polynomial, or is actually polynomial. Since our methods work just as well in the more general setting of continuous mixtures, we believe that doing better than quasi-polynomial complexity would require an algorithm specific to a mixture of $k$ Gaussians. Further structure, e.g., hierarchical structure, could also make the problem easier, though our current analysis does not benefit from such assumptions.

We note that mixtures of Gaussians are particularly suited to learning with diffusion models because diffusions preserve the Gaussian. We expect that similar results are possible in other settings with a "match" between the family and the diffusion, e.g., mixtures of product distributions on the hypercube where diffusion is a random walk on the hypercube (i.e. bit-flip noise). It would be interesting to find other families of distributions which be learned using diffusion models, including families of conditional distributions such as mixtures of linear regressions.

Finally, we note that in contrast to our algorithm based on piecewise polynomial regression, in practice the score function is typically learned with a neural network. The score function for a Gaussian mixture is exactly represented by a softmax neural network, which raises the question of whether it can be learned by a neural network with gradient descent. Understanding neural network training dynamics for diffusion models is an important open direction.

## 1.1 Related work

**Learning mixtures of Gaussians.** The problem of learning a mixture of Gaussians from samples has an illustrious history [TSM85]. We first distinguish between several types of results. First, one can ask for parameter learning or distribution learning—either outputting parameters that are close to the ground-truth parameters, or simply a distribution that is close (e.g., in TV distance) to the ground-truth distribution. For distribution learning, the learning can be improper, that is,

the output need not be a mixture of Gaussians. Second, one can either study the problem from an information theoretic perspective—the minimum number of samples required to get within a specified error regardless of computational cost—or computational complexity perspective—where the emphasis is on the running time of the algorithm. Below, $n$ denotes the ambient dimension, $k$ the number of components, and $\varepsilon$ the target accuracy.

Most earlier works go through parameter learning and require the means to be sufficiently separated. For identity-covariance Gaussians, [DS00, AK05, VW04] show that spectral methods work with a separation of at least $\widetilde{\Omega}(\min\{n, k\}^{1/4})$. Allowing arbitrary covariances, [MV10] show that whenever the mixture is "$\varepsilon$-statistically learnable," it can be learned with running time and sample complexity $\exp(k) \operatorname{poly}(n, \frac{1}{\varepsilon})$. Several works use the sum-of-squares method [KSS18, HL18] to learn a mixture with separation $\Omega(k^\gamma)$ in time $n^{\operatorname{poly}(1/\gamma)}$. These methods extend to a broader class of mixture distributions, where components have moments which can be certifiably bounded. [LL22] obtain a polynomial-time algorithm whenever the separation is $\Omega(\ln^{\frac{1}{2}+c} k)$ for constant $c > 0$.

Separation conditions are unavoidable for sample-efficient parameter learning. [RV17] shows that the threshold for efficient parameter learning is $\Theta(\sqrt{\ln k})$: with separation $\Omega(\sqrt{\ln k})$, poly-nomially many samples suffice (information-theoretically), while with separation $o(\sqrt{\ln k})$, super-polynomially many samples are required. [DWYZ20] conducts a more fine-grained study of this problem.

Hence, any sample-efficient algorithm for learning mixtures of Gaussians without separation cannot go through parameter learning. The optimal information theoretic complexity is known up to logarithmic factors: [ABDH$^+$18] show a sample complexity bound in TV distance of $\widetilde{\Theta}(kn^2/\varepsilon^2)$ for a mixture of $k$ Gaussians in $\mathbb{R}^n$ (with any variance) and $\widetilde{\Theta}(kn/\varepsilon^2)$ for axis-aligned Gaussians. However, their algorithms are based on brute-force search and have exponential running time. [ADLS17] give an improper nearly linear-time algorithm based on polynomial interpolation which learns a mixture of Gaussians with arbitrary variances in 1 dimension, with $\widetilde{O}(k/\varepsilon^2)$ samples. Most relevant for us, for a mixture of Gaussians with identity covariance, a breakthrough work by [DK20] uses algebraic geometry to obtains a time and sample complexity of $(k/\varepsilon)^{O(\ln^2 k)}$ plus polynomial factors.

In the statistics literature, the model is referred to as the Gaussian location mixture. [SG20, KG22] consider arbitrary mixing measures $Q_0$ and give finite-sample bounds using non-parametric MLE for squared Hellinger risk, which scale as $\frac{\log N}{N}$ in terms of the number of samples $N$. The risk depends on the volume of an approximate support of $Q_0$, but also include a constant with unspecified dependence on the dimension $n$.

We note that the picture is more complicated when variances are unknown: [DKS17] obtain statistical query (SQ) lower bounds ($2^{n^{\Omega(1)}}$ queries of fixed polynomial precision) based on a connection with non-Gaussian component analysis, with a "parallel pancake" construction in a unknown direction. [GVV22] show that $\log n$ components are enough to obtain a super-polynomial lower bound assuming exponential hardness of the classical LWE problem.

Finally, a recent line of work considers the problem of robust learning of Gaussian mixtures [LM21, BDJ$^+$22], i.e., under adversarial corruption of some fraction of samples; these methods have complexity $n^{O(k)}$.

**Concurrent work.** During the preparation of this manuscript, we were made aware of independent and concurrent work by Chen, Kontonis, and Shah [CKS24] which also gave guarantees for learning Gaussian mixtures using diffusion models and piecewise polynomial regression for score es-

timation. They consider the more general case where covariances are well-conditioned but arbitrary, but their runtime scales exponentially in $\text{poly}(k/\epsilon)$ rather than $\text{poly}\log(k/\epsilon)$, which is unavoidable by SQ lower bounds [DKS17]. In contrast to our work, their work uses a different approach to polynomial approximation and does one-shot learning of parameters via spectral methods.

**Diffusion models.** Diffusion models [SDWMG15, SE19, SSDK+20] are a modern generative modeling paradigm which involves defining a forward noising process which turns a data distribution into a pure noise (e.g., Gaussian) distribution, and then learning to simulate the reverse process. For SDE-based diffusion models, the reverse process involves the score function (gradient of log-pdf) of the intermediate distributions; hence they are also called score-based generative models (SGM). See [TZ24] for a technical tutorial. We note that diffusion models are essentially a reparameterization of stochastic localization [Eld13] as pointed out by [Mon23]. Stochastic localization has been independently studied in the probability literature and been used to obtain new results in sampling [CE22, EAMS22].

Theoretical work has focused on two problems: (1) When is it possible to efficiently learn the score? (2) Given a learned ($L^2$-accurate) score function, what guarantees can we obtain for sampling from the data distribution? Answers to these two questions together would give an end-to-end result for learning via diffusion models.

Addressing (2), it is a remarkable fact that having a $L^2$-accurate score function for the sequence of distributions is sufficient for sampling under minimal distributional assumptions, allowing even multimodal distributions which cause slow mixing for local MCMC algorithms [LLT23, CHZW23, CLL23]. [BDBDD23] show that it suffices to have a number of steps linear in the dimension.

Question (1) has proved to be thornier; it has been a challenge to obtain end-to-end results for non-trivial settings. Several works consider the problem of representability by neural networks, such as [CL24] for distributions whose log-density relative to a Gaussian can be represented by a low-complexity neural network, or [MW23] for graphical models. Following the work on neural network function approximation for smooth functions, [OAS23] give nearly minimax optimal estimation rates for densities in Besov spaces. [WWY24] relates score learning to kernel density estimation. The manifold assumption is another popular setting for analysis: [DB22] considers generalization error, and [CHZW23] give learning guarantees when the distribution is supported on a subspace (however, this family of distributions can be trivially learned by first recovering the subspace). Finally, we note that score matching can be used to learn exponential families for which sampling is difficult [KHR22, PRS+24]; however, these methods only use the score for the data distribution, rather than a sequence of distributions as in a diffusion model.

[SCK23] consider using diffusion models to learn Gaussian mixtures, and show that diffusion models can do as well as the EM algorithm. However, they either require $k = 2$, or the components to be well-separated and $O(1)$-warm starts to be given for all the means. Gaussian mixtures are a popular toy model for understanding various aspects or behavior of diffusion models, including learning behavior [CKVEZ23], sample complexity [BM23], guided diffusion [WCL+24], and critical windows [LC24].

**Analytic conditions for learning functions.** We rely on the noise sensitivity/stability framework [KOS08], which shows that Gaussian noise stability (or small Gaussian surface area) implies approximability by a low-degree polynomial, giving an efficient "low-degree algorithm" for learning. We note a similar-in-spirit result that under the Gaussian distribution, intersections of $k$

halfspaces can be learned in time $n^{O(\ln k)}$. The noise sensitivity framework was previously developed for boolean functions on the hypercube [BKS99] and applied to learning function classes such as functions of halfspaces [KOS04, KKMS08].

Learnability by neural networks can also be related to complex analytic properties of the function: [SZ23] relate a radius of analyticity to the Hermite expansion which gives results on representability by neural networks. Learnability by neural networks for multi-index models [BBPV23]—functions depending on the projection of the input to a few dimensions (such as the score function of a Gaussian mixture with $k \ll n$)—is also related to the Hermite expansion of the function.

## 1.2 Notation

We let $\gamma_{\mu,\sigma^2}$ denote the density of $\mathcal{N}(\mu, \sigma^2 I_n)$, and abbreviate $\gamma_{\sigma^2} = \gamma_{0,\sigma^2}$. We abbreviate this as $\gamma$ when $\sigma$ is understood. In general, we denote probability measures by uppercase letters and their corresponding densities by lowercase letters, though sometimes we will conflate the two.

We use $|v| = |v|_2$ to denote the norm of a vector $v \in \mathbb{R}^n$, to avoid confusion with function norms. For a measure $\nu$ on $\Omega$, let $\|f\|_{L^p(\nu)} = \left( \int_\Omega f^p \, d\nu \right)^{1/p}$. When the measure is clear, we may simple write $\|f\|_p$. Let $\|f\|_\nu := \|f\|_{L^2(\nu)}$. For a $\mathbb{R}^n$-valued function, we write $\|f\|_{L^p(\nu)}$ to mean $\||f|\|_{L^p(\nu)}$. For $x_0 \in \mathbb{R}^n$, let $B_R(x_0) = \{x \in \mathbb{R}^n : |x - x_0| \leq R\}$, and let $B_R = B_R(0)$.

Let $\binom{S}{k}$ denote the set of subsets of $S$ of size $k$, and $\binom{S}{\leq k}$ denote the set of subsets of $S$ of size at most $k$.

We give some further background and notation on Markov semigroups and generators in Section 4.

## 2 Proof overview

The idea behind diffusion models is to first define a forward process based on a SDE that takes the data distribution to a pure noise distribution, in our case a Gaussian. Then using a result on reversing a SDE [And82], we can write down the reverse process which involves the score function. We apply this in the case when the forward process is simply Brownian motion,

$$dx_t = dW_t, \quad x_0 \sim P_0$$

to obtain that this process on $[0, T]$ is equivalently described by

$$dx_t = \nabla \ln p_t(x_t) + d\widetilde{W}_t, \quad 0 \leq t \leq T, \quad x_T \sim P_T$$

where $x_t$ has distribution $p_t$ and $\widetilde{W}_t$ is reverse Brownian motion. We choose $T$ large enough so that $P_T$ is close to Gaussian. Hence, if we learn the score functions $\nabla \ln p_t(x_t)$, then we can approximately simulate the reverse process. In Section 3, we make this precise by adapting known convergence results on diffusion models, which show that we can approximately sample from the data distribution given a $L^2$-accurate score function.

By Tweedie's formula, the score function admits an interpretation in terms of posterior mean given an observation, so the score matching objective $\mathbb{E}_{p_t} |\nabla \ln p_t - s|^2$ can be rewritten as the *supervised* denoising auto-encoder objective (see Section 3.2). It now remains to give a way to learn the score function, by showing that it is in a low-dimensional function class we can optimize over.

In Section 4, we show that in the special case where all centers are close together, there is a polynomial of low (poly log $k$) degree that approximates the score function with respect to the mixture distribution. The key enabling result is Lemma 4.1, which shows that Gaussian noise stability implies low-degree approximability, by considering the Hermite expansion. In contrast to existing literature, we use a higher-order version of noise stability which involves bounding $\mathscr{L}^m f$ where $\mathscr{L}$ is the generator of the Ornstein-Uhlenbeck process. The interpretation of the score function as a posterior mean gives us a handle on computing (Lemma 4.2 in Section 4.1) and bounding (Lemma 4.4 in Section 4.2) its derivatives. However, one key problem remains: Lemma 4.4 gives us low-degree approximability with respect to the Gaussian, not the actual mixture $P_t$. A simple change-of-measure argument fails to work because because the degree of the polynomial required depends poly-logarithmically on the desired accuracy $\varepsilon$, while the change of measure gives a multiplicative factor exponential in the degree, resulting in a factor $\varepsilon e^{\ln^C(1/\varepsilon)} \gg 1$. Instead, we have to smooth the distribution by the Ornstein-Uhlenbeck process *before* the polynomial approximation; moreover, we need to consider a higher-order smoothing (Lemma 4.5 in Section 4.3).

In Section 5 we consider the general setting of multiple clusters. Imagine running the backwards diffusion process from pure Gaussian noise back to our mixture: the idea is to inductively maintain "warm starts" to all the Gaussian centers at the current resolution, and perform piecewise polynomial regression in Voronoi cells centered at these warm starts. If we have the warm starts, then we can approximate the score function with "local" score functions on Voronoi cells (Lemma 5.3 in Section 5.1) and piggybacking off the results of Section 4, obtain a low-degree piecewise polynomial approximation (Lemma 5.6 in Section 5.2). We analyze the sample complexity of piecewise polynomial regression in Section 5.3 (Lemma 5.8) and address the problem of maintaining warm starts in Section 5.4 (Lemma 5.13); this we get inductively from the score estimates as they exactly point in the "denoising" direction, i.e., approximately towards these centers. With all these pieces, we can then prove the main theorem.

*Proof of Theorem 1.2.* Here we show that Algorithm 1 successfully samples from $P_0 = Q_0 * \mathcal{N}(0, \sigma_0^2)$ with the claimed time and sample complexity. Let $M_2 = \mathbb{E}_{x \sim P_0} |x|^2$. Let $\widehat{P}_{t_1}$ be the output of Algorithm 1. By the triangle and Pinsker's inequality,

$$\mathrm{TV}(P_0, \widehat{P}_{t_1}) \le \mathrm{TV}(P_0, P_{t_1}) + \mathrm{TV}(P_{t_1}, \widehat{P}_{t_1}) = \mathrm{TV}(P_0, P_{t_1}) + \sqrt{\frac{1}{2} \mathrm{KL}(P_{t_1} \| \widehat{P}_{t_1})}.$$

Choose the starting time $t_1 = \frac{\varepsilon^2 \sigma_0^2}{2\sqrt{n}}$, so that by Lemma 3.3 we have $\mathrm{TV}(P_0, P_{t_1}) \le \frac{\varepsilon^2}{2}$. Hence it is sufficient to prove $\mathrm{KL}(P_{t_1} \| \widehat{P}_{t_1}) \le \frac{\varepsilon^2}{2}$; actually it suffices to find parameters to make $\mathrm{KL}(P_{t_1} \| \widehat{P}_{t_1}) = O(\varepsilon^2)$ as the exact constant can be adjusted by rescaling $\varepsilon$ by a constant. We implement Algorithm 1 with the step size schedule obtained recursively by the equality version of condition 3 in Theorem 3.1, i.e. $t_k + 1 = (t_{k+1} + 1) \max\{e^{-2\kappa}, (t_{k+1} + 1)^{-\kappa}\}$ for $\kappa = \frac{\varepsilon^2}{M_2 + n \ln(T+1)}$, with ending time $t_{N_{\mathrm{step}}} = T = \frac{M_2 + d}{\varepsilon^2}$, and number of iterations $N_{\mathrm{step}} = O\left(\frac{1}{\kappa} \ln\left(\frac{T+1}{t_1}\right)\right)$. Then, given that we pick $\varepsilon_\ell^2 = \frac{\varepsilon^2}{\ln(T+1)}$, Theorem 3.1 tells us that $\mathrm{KL}(P_{t_1} \| \widehat{P}_{t_1}) = O(\varepsilon^2)$, as needed.

Hence the problem reduces to giving sufficiently accurate estimates of the score function, to guarantee an accuracy of $\varepsilon_\ell^2 = \frac{\varepsilon^2(\sigma_\ell^2 + 1)}{\ln(T+1)}$ for all time steps $t_1, \ldots, t_{N_{\mathrm{step}}}$. Obtaining such accurate scores is our main contribution and the proof consists of showing by backwards induction on $\ell = N_{\mathrm{step}}, \ldots, 1$ that the following hold:

9

1. The set $\mathcal{C}_\ell$ is a complete set of $R_\ell$-warm starts for $Q_0$ (see Definition 5.1), where $R_\ell :=$ $C\left(R_0 + 2\sigma_\ell\sqrt{\ln\left(\left(\frac{R_0}{\sigma_\ell}+1\right)\frac{k}{\alpha_{\min}}\right)}\right)$ for a large enough universal constant $C$. Moreover, $|\mathcal{C}_\ell| \leq$ $k' = O\left(k\ln\left(\frac{1}{\alpha_{\min}}\right)\right)$.

2. The score estimate $s_{t_\ell}$ is $\varepsilon_\ell$-accurate in $L^2(P_{t_\ell})$, for $\varepsilon_\ell^2 = \frac{\varepsilon^2(\sigma_\ell^2+1)}{\ln(T+1)}$.

More precisely, we will show that 1 holds for $\ell = N_{\text{step}}, \ldots, N_{\text{step}} - \ell'$ and 2 holds for $\ell = N_{\text{step}}, \ldots, N_{\text{step}} - \ell' + 1$ with probability at least $1 - \frac{\ell'\delta}{N_{\text{step}}}$.

The base case of the induction ($\ell = N_{\text{step}}$, $\ell' = 0$) is to show that $\mathcal{C}_{N_{\text{step}}} = \{0\}$ is a complete set of $R_{N_{\text{step}}}$-warm-starts. But note that the assumption $\varepsilon^2 \leq \frac{M_2+n}{D^2}$ implies the variance at step $N_{\text{step}}$ satisfies $\sigma_{N_{\text{step}}}^2 = T + \sigma_0^2 \geq T = \frac{M_2+n}{\varepsilon^2} \geq D^2$, which given that we pick $C \geq 1$ means $R_{N_{\text{step}}} \geq \sigma_{N_{\text{step}}} \geq D$. Therefore, from the definition of $D$, we have that $Q_0$ is supported on $B_D \subseteq B_{R_{N_{\text{step}}}}$.

Next, we show the induction step; here, the hypothesis of induction for step $\ell$ is that the set $\mathcal{C}_\ell$ is a complete set of $R_\ell$-warm starts for $Q_0$. Then, we show that 2 for $\ell$ and 1 for $\ell-1$ are satisfied, after excluding an event of probability at most $\frac{\delta}{N_{\text{step}}}$. First we handle 2. Using Lemma 5.6 with

$$\widetilde{\varepsilon}^2 := \frac{\varepsilon^2}{\left(\left(\frac{R_\ell}{\sigma_\ell}\right)^2 + \ln\left(\frac{k'R_\ell(M_2+n)}{\sigma_\ell\varepsilon}\right)\right)\ln(T+1)},$$

there exists a piece-wise polynomial $\widetilde{g}_\ell$ on the Voronoi partition of $\mathcal{C}_\ell$ that approximates the score function with the desired accuracy,

$$\left\|\widetilde{g}_\ell - f_{\sigma_\ell^2}\right\|_{P_{t_\ell}}^2 \lesssim \widetilde{\varepsilon}^2\left(R_\ell^2 + \sigma_\ell^2\ln\left(\frac{k'}{\widetilde{\varepsilon}}\right)\right) \lesssim \frac{\varepsilon^2(\sigma_\ell^2+1)}{\ln(T+1)}. \tag{1}$$

We note

$$\ln\left(\frac{1}{\widetilde{\varepsilon}}\right) = O\left(\ln\frac{\left(\frac{R_\ell}{\sigma_\ell}\right)(M_2+n)}{\varepsilon}\right) = O\left(\ln\left(\frac{1}{\varepsilon}\right)\right)$$

assuming that $\varepsilon \leq \min\left\{\frac{1}{2}, \frac{\sigma_0}{R_0}, \frac{1}{D}, \frac{1}{n}, \alpha_{\min}\right\}$ and noting $\sigma_0 \leq \sigma_\ell$, $\alpha_{\min} \leq \frac{1}{k}$. The degree of $\widetilde{g}_\ell$ is at most

$$d_\ell = O\left(\left(\frac{R_\ell}{\sigma_\ell} + \sqrt{\ln\left(\frac{k'}{\widetilde{\varepsilon}}\right)}\right)^6 \ln\left(\frac{1}{\widetilde{\varepsilon}}\right)^4\right)$$

$$= O\left(\left(\frac{R_0}{\sigma_\ell} + \sqrt{\ln\left(\left(1 + \frac{R_0}{\sigma_\ell}\right)\frac{k}{\alpha_{\min}}\right)} + \sqrt{\ln\left(\frac{1}{\varepsilon}\right)}\right)^6 \ln\left(\frac{1}{\varepsilon}\right)^4\right)$$

$$= O\left(\left(\ln\left(\frac{1}{\varepsilon}\right)^3 + \left(\frac{R_0}{\sigma_0}\right)^6\right)\ln\left(\frac{1}{\varepsilon}\right)^4\right). \tag{2}$$

10

Now for an arbitrary warm start point $\widehat{\mu}_i^{(\ell)} \in \mathcal{C}_\ell$, by Lemma 5.6, there exists a polynomial $\widetilde{g}_\ell = \sum_{|\mathbf{k}| \leq d_\ell} \widetilde{b}_\mathbf{k}^{(i)} h_\mathbf{k}(y - \widehat{\mu}_i^{(\ell)})$ where $h_\mathbf{k} = h_{\mathbf{k}, \sigma_\ell^2}$ are the Hermite polynomials with variance $\sigma^2$, and the coefficients satisfy

$$\left| (\widetilde{b}_\mathbf{k}^{(i)})_{|\mathbf{k}| \leq d_\ell} \right| = \|\widetilde{g}_\ell\|_{\gamma_{\widehat{\mu}_i, \sigma_\ell^2}} \leq D.$$

But letting $d_\ell$ be the RHS of (2) with appropriate constants, the condition $d_\ell \geq (R_\ell/\sigma_\ell)^2$ is satisfied, and hence we can use Lemma 5.8 with parameters $D, M$ both set to $D$; the implication is that given

$$N = \varepsilon'^{-4} n \left( k(D+1)^2 \ln\left(\frac{1}{\alpha_{\min}}\right) \right)^2 \left( n \ln\left(\frac{k' N_{\text{step}}}{\delta}\right) \right)^{cd_\ell}$$

samples for $\varepsilon'^2 = \frac{\varepsilon^2 (\sigma_\ell + 1)}{\ln(T+1)}$ and some universal constant $c$, we have the guarantee

$$\|\widehat{g}_\ell - f\|_{L^2(P_{t_\ell})}^2 - \|\widetilde{g}_\ell - f\|_{L^2(P_{t_\ell})}^2 \leq \varepsilon'^2$$

after excluding an event of probability at most $\frac{\delta}{2N_{\text{step}}}$. Combining this with Equation (1), we obtain

$$\|\widehat{g}_\ell - f\|_{L^2(P_{t_\ell})}^2 \leq O(\varepsilon'^2)$$

as desired.

Next, we handle 1, i.e. show that $\mathcal{C}_{\ell-1}$ is a complete set of $R_{\ell-1}$-warm starts for $Q_0$. Note that when $\sigma_\ell$ changes by a constant factor, 1 is still satisfied with a modified constant, so we only have to update $\mathcal{C}_\ell$ each time $t_\ell + 1$ halves (as done in Section 1). According to Lemma 5.13, it is sufficient to find an estimate $g$ for the score which satisfies $\|f - g\|_{L^2(P_{t_\ell})}^2 \leq (R_0 + \sigma_\ell)^2 \alpha_{\min}$. But from our assumption we have $\varepsilon'^2 = \frac{\varepsilon^2 (\sigma_\ell + 1)}{\ln(T+1)} \leq (R_0 + \sigma_\ell)^2 \alpha_{\min}$, which means the polynomial $\widetilde{g}_\ell$ that we obtained for proving part 2 in the induction already satisfies the required accuracy for updating the warm starts as well; hence we can use the same degree of polynomial for our regression here. Again we exclude an event of probability at most $\frac{\delta}{2N_{\text{step}}}$. Note that the sample size required by Algorithm 2 is $O\left(\frac{\ln(N_{\text{step}}/\delta)k}{\alpha_{\min}}\right)$, which is negligible compared to the number of samples needed for regression. This completes the induction step.

Multiplying $N$ by the number of steps $N_{\text{step}}$ and dropping lower-order terms (recalling the assumption on $\varepsilon$), we get that the total sample complexity is

$$N^{\text{total}} = \left( n \ln\left(\frac{1}{\delta}\right) \right)^{O(d_\ell)} = \left( n \ln\left(\frac{1}{\delta}\right) \right)^{O\left( \left( \ln\left(\frac{1}{\varepsilon}\right)^3 + \left(\frac{R_0}{\sigma_0}\right)^6 \right) \ln\left(\frac{1}{\varepsilon}\right)^4 \right)}.$$

$\square$

Note that the regression problems that we solve to obtain score estimates for running DDPM can use the same batch of samples, because we can still do a union bound. However, we do need to use fresh samples each time for updating the warm starts, separate from the ones that we use for score estimates. This is because our uniform convergence bounds in Lemma 5.8 require the condition that the samples are independent from the randomness used in the construction of the Voronoi partition.

*Proof of Corollary 1.3.* This follows from Theorem 1.2 with $\alpha_{\min} = \frac{\varepsilon}{C^d}$ and substituting $\min\left\{\frac{\varepsilon}{C^d}, \frac{1}{2}, \frac{1}{R_0}, \frac{1}{n}, \frac{1}{D}\right\}$ for $\varepsilon$ in Theorem 1.2, noting $\ln\left(\frac{C^d}{\varepsilon}\right) = O\left(d + \ln\left(\frac{1}{\varepsilon}\right)\right)$. □

# 3  Diffusion models

Using the framework of diffusion models, the algorithm and analysis consist of two main parts: showing that we can learn this score function efficiently, and showing that an accurate score function allows generating more samples from the distribution. In this section, we focus on the second part, which has been well-studied.

## 3.1  Convergence guarantees for diffusion models

Two popular parameterizations of diffusion models are the variance-exploding (VE) and variance-exploding (VP) processes. The diffusion models are also known as score matching with Langevin dynamics (SMLD) and denoising diffusion probabilistic models (DDPM), respectively. The forward SDE for the VE process is simply Brownian motion, while the forward SDE for the VP process is the Ornstein-Uhlenbeck process. Below, we lay out the forward SDE for each of these process, and the backward SDE obtained by [And82]:

VE:

$$dy_t = dW_t$$
$$dy_t = -\nabla \ln p_t(y_t)\, dt + d\widetilde{W}_t$$
$$dy_t^{\leftarrow} = \nabla \ln p_{T-t}(y_t^{\leftarrow})\, dt + dW_t'$$

VP:

$$dx_t = -x_t\, dt + \sqrt{2}\, dW_t$$
$$dx_t = (-x_t - 2\nabla \ln p_t(x_t))\, dt + \sqrt{2}\, d\widetilde{W}_t$$
$$dx_t^{\leftarrow} = (x_t^{\leftarrow} + 2\nabla \ln p_{T-t}(x_t^{\leftarrow}))\, dt + \sqrt{2}\, dW_t'$$

where $p_t$ is the density of the either $x_t$ or $y_t$, $\widetilde{W}_t$ is reverse Brownian motion, and $W_t$, $W_t'$ are usual Browian motions, and the reverse processes are initialized at $y_0^{\leftarrow} \sim p_T$, $x_0^{\leftarrow} \sim p_T$ and we can match up trajectories $y_t^{\leftarrow} = y_{T-t}$, $x_t^{\leftarrow} = x_{T-t}$ for $t \in [0, T]$. For convenience, we work with the VE process. As most convergence guarantees in the literature are for the VP process, we need to adapt those results. Note that as continuous processes we have $x_t = e^{-t} y_{e^{2t}-1}$, but some care is required when matching up the discretizations.

The following is an adaptation of [BDBDD23, Theorem 1], after reparameterizing the VP into the VE process.

**Theorem 3.1** (Reverse KL guarantee for variance-exploding diffusion models). *Let $0 < t_1 < t_2 < \ldots < t_{N_{\text{step}}} = T$ and $h_k = t_{k+1} - t_k$. Suppose $T \geq 1$. Assume the following.*

1. *We have a score function estimate $s_t(y)$ for each $t = t_k$ such that*

$$\mathbb{E}_{P_{t_k}} |\nabla \ln p_{t_k}(y) - s_{t_k}(y)|^2 \leq \varepsilon_k^2.$$

2. *The data distribution has bounded second moment $M_2 = \mathbb{E}_{P_0} |y|^2$.*

3. *For some $\kappa < 1$, the step size schedule satisfies*

$$t_k + 1 \geq (t_{k+1} + 1) \max\left\{e^{-2\kappa}, (t_{k+1} + 1)^{-\kappa}\right\}.$$

Let $\widehat{p}_t$ denote the distribution of the following discretization of the reverse process when initialized at $\widehat{p}_T = \mathcal{N}(0, (T+1) \cdot I_n)$:[1]

$$\widehat{y}_{t_k} = \widehat{y}_{t_{k+1}} + 2\left[(t_{k+1} - 1) - \sqrt{(t_k - 1)(t_{k+1} - 1)}\right] s_{t_{k+1}}(\widehat{y}_{t_{k+1}}) + \sqrt{h_k} \cdot \xi_{t_k}.$$

Then

$$\mathrm{KL}(p_{t_1} \| \widehat{p}_{t_1}) \lesssim \frac{n + M_2}{T+1} + \sum_{k=1}^{N_{\text{step}}-1} \ln\left(\frac{t_{k+1}+1}{t_k+1}\right) \cdot \frac{1}{t_k+1} \cdot \varepsilon_{k+1}^2 + \kappa n \ln(T+1) + \kappa^2 n N_{\text{step}} + \kappa M_2.$$

(3)

Moreover, we can choose a schedule of length $N_{\text{step}} = O\left(\frac{1}{\kappa} \ln\left(\frac{T+1}{t_1}\right)\right)$ to make assumption 3 hold. If we choose $T = \frac{M_2+n}{\varepsilon^2}$ and $\kappa = \frac{\varepsilon^2}{M_2+n \ln(T+1)}$, and we have $\varepsilon_k^2 \leq \frac{\varepsilon^2(t_k+1)}{\ln(T+1)}$ for each $k$, then the error is $O(\varepsilon^2)$.

We conclude this from the analogous guarantee for the VP process.

**Theorem 3.2** (Reverse KL guarantee for variance-preserving diffusion models [BDBDD23, Theorem 2]). *Let* $0 < t_1 < t_2 < \ldots < t_{N_{\text{step}}} = T$ *and* $h_k = t_{k+1} - t_k$. *Suppose* $T \geq 1$. *Let* $P_t^{\text{OU}}$ *denote the distribution of the Ornstein-Uhlenbeck process run for time $t$ initialized from $P_0$, that is, the distribution of* $e^{-t}X + \sqrt{1 - e^{-2t}} \cdot \xi$ *where* $X \sim P_0$ *and* $\xi \sim \mathcal{N}(0, I_d)$. *Assume the following.*

1. *We have a score function estimate $s_t(x)$ for each $t = t_k$ such that*

$$\mathbb{E}_{P_{t_k}^{\text{OU}}} \left| \nabla \ln p_{t_k}^{\text{OU}}(x) - s_{t_k}(x) \right|^2 \leq \varepsilon_k^2.$$

2. *The data distribution has bounded second moment* $M_2 = \mathbb{E}_{P_0} |x|^2$.

3. *The step size schedule satisfies* $h_k \leq \kappa \min\{1, t_k\}$ *(i.e.* $h_k \leq \min\left\{\kappa, \frac{\kappa}{\kappa+1} t_{k+1}\right\}$*)*

Let $\widehat{p}_t$ denote the distribution of the exponential integrator discretization of the reverse process when initialized at $\widehat{P}_T = \mathcal{N}(0, I_n)$, given by

$$\widehat{x}_{t_k} = e^{h_k}\widehat{x}_{t_{k+1}} + 2(e^{h_k} - 1)s_{t_{k+1}}^x(\widehat{x}_{t_{k+1}}) + \sqrt{e^{2h_k} - 1} \cdot \xi_{t_k}, \quad \xi_{t_k} \sim \mathcal{N}(0, I_n).$$

Then

$$\mathrm{KL}(p_{t_1} \| \widehat{p}_{t_1}) \lesssim (n + M_2)e^{-2T} + \sum_{k=1}^{N_{\text{step}}} h_k \varepsilon_k^2 + \kappa n T + \kappa^2 n N_{\text{step}} + \kappa M_2.$$

(4)

Note that as opposed to [BDBDD23], we index time in the forward direction.

---

[1]It would be more natural to have $h_k$ as the coefficient in front of the score term. This coefficient comes from doing a change-of-variable from Theorem 3.2 for the DDPM process, rather than reproving the theorem for this process.

*Proof of Theorem 3.1.* Consider the following two sets of continuous processes,

$$dx_t = -x_t \, dt + \sqrt{2} \, dW_t \qquad\qquad\qquad dy_t = dW_t$$

with same initial distribution $x_0, y_0 \sim P$. Let $P_t^x$ and $P_t^y$ be the distribution of the two processes at time $t$. We claim that

$$x_t = e^{-t} y_{e^{2t}-1}.$$

To see this, we check that $e^{-t} y_{e^{2t}-1}$ satisfies the equation for $x_t$:

$$
\begin{aligned}
d(e^{-t} y_{e^{2t}-1}) &= -e^{-t} y_{e^{2t}-1} \, dt + e^{-t} \, dy_{e^{2t}-1} \\
&= -(e^{-t} y_{e^{2t}-1}) \, dt + e^{-t} \sqrt{2e^{2t}} \, dW_t \\
&= -(e^{-t} y_{e^{2t}-1}) \, dt + \sqrt{2} \, dW_t.
\end{aligned}
$$

By the change-of-variables formula,

$$
\begin{aligned}
p_t^x(x) &= e^{tn} p_{e^{2t}-1}(e^t x) \\
\implies \nabla \ln p_t^x(x) &= e^t \nabla \ln p_{e^{2t}-1}(e^t x).
\end{aligned}
$$

Write $s_t^y = s_t$ for clarity. Defining $s_t^x$ to satisfy

$$s_t^x(x) = e^t s_{e^{2t}-1}^y(e^t x),$$

we have

$$\mathbb{E}_{P_t^x} |\nabla \ln p_t^x(x) - s_t^x(x)|^2 = e^{2t} \mathbb{E}_{P_t^y} \left| \nabla \ln p_{e^{2t}}^y(y) - s_{e^{2t}-1}^y(y) \right|^2.$$

Now define

$$t_k^x = \frac{\ln(t_k + 1)}{2}, \qquad\qquad h_k^x = t_{k+1}^x - t_k^x,$$

and consider the discrete processes defined backwards in time,

$$
\begin{aligned}
\widehat{x}_t &= e^h \widehat{x}_{t+h} + 2(e^h - 1) s_{t+h}^x(\widehat{x}_{t+h}) + \sqrt{e^{2h} - 1} \cdot \xi_t^x & (t, h) = (t_k^x, h_k^x) \\
\widehat{y}_t &= \widehat{y}_{t+h} + 2\left[(t+h-1) - \sqrt{(t-1)(t+h-1)}\right] s_{t+h}^y(\widehat{y}_{t+h}) + \sqrt{h} \cdot \xi_t & (t, h) = (t_k, h_k),
\end{aligned}
$$

where $x_{t_{N_{\text{step}}}^x} \sim P_{t_{N_{\text{step}}}^x}^x$ and $y_{t_{N_{\text{step}}}} = e^{t_{N_{\text{step}}}^x} x_{t_{N_{\text{step}}}^x}$ (so that $y_{t_{N_{\text{step}}}} \sim P_{t_{N_{\text{step}}}}^y$), and we couple $\xi_{t_k^x}^x = \xi_{t_k}$. We can inductively check that

$$\widehat{x}_{t_k} = e^{-t_k^x} \widehat{y}_{t_k}.$$

Indeed, if this holds for $k + 1$, then

$$
\begin{aligned}
\widehat{x}_{t_k^x} &= e^{h_k^x} \widehat{x}_{t_{k+1}^x} + 2(e^{h_k^x} - 1) s_{t_{k+1}^x}(\widehat{x}_{t_{k+1}^x}) + \sqrt{e^{2h_k^x} - 1} \cdot \xi_t^x \\
&= e^{h_k^x}(e^{-t_{k+1}^x} \widehat{y}_{t_{k+1}}) + 2(e^{h_k^x} - 1) e^{t_{k+1}^x} s_{t_{k+1}}^y(\widehat{y}_{t_{k+1}}) + \sqrt{e^{2h_k^x} - 1} \cdot \xi_t^x \\
&= e^{-t_k^x} \left( \widehat{y}_{t_{k+1}} + 2(e^{h_k^x} - 1) e^{t_{k+1}^x + t_k^x} s_{t_{k+1}}^y(\widehat{y}_{t_{k+1}}) + \sqrt{t_{k+1} - t_k} \cdot \xi_t \right) \\
&= e^{-t_k^x} \left( \widehat{y}_{t_{k+1}} + 2\left[(t_{k+1} - 1) - \sqrt{(t_k - 1)(t_{k+1} - 1)}\right] s_{t_{k+1}}^y(\widehat{y}_{t_{k+1}}) + \sqrt{t_{k+1} - t_k} \cdot \xi_t \right) \\
&= e^{-t_k^x} \widehat{y}_{t_k}.
\end{aligned}
$$

14

Hence, by Theorem 3.2,

$$\mathrm{KL}(p^y_{t_1}\|\widehat{p}^y_{t_1}) = \mathrm{KL}(p^x_{t_1}\|\widehat{p}^x_{t_1}) \lesssim (n+M_2)e^{-2t^x_{N_{\text{step}}}} + \sum_{k=1}^{N_{\text{step}}} h^x_k \cdot e^{-2t^x_k} \cdot \varepsilon^2_k + \kappa n t^x_{N_{\text{step}}} + \kappa^2 n N_{\text{step}} + \kappa M_2$$

$$\lesssim \frac{n+M_2}{T+1} + \sum_{k=1}^{N_{\text{step}}} \ln\left(\frac{t^y_{k+1}+1}{t^y_k+1}\right) \cdot \frac{1}{t^y_k+1} \cdot \varepsilon^2_k + \kappa n \ln(T+1) + \kappa^2 n N_{\text{step}} + \kappa M_2.$$

Because $y_T = \sqrt{T+1}\cdot x_{\frac{\ln(T+1)}{2}}$, the initialization $x_{t^x_N} \sim \mathcal{N}(0,I)$ corresponds to $y_T \sim \mathcal{N}(0,(T+1)\cdot I)$. The requirement on step sizes is

$$\frac{\ln(t_{k+1}+1)-\ln(t_k+1)}{2} \le \min\left\{\kappa, \frac{\kappa}{\kappa+1}\cdot\frac{\ln(t_{k+1}+1)}{2}\right\}$$

$$\iff t_{k+1}+1 \ge (t_{k+1}+1)\max\{e^{-2\kappa}, (t_{k+1}+1)^{-\frac{\kappa}{\kappa+1}}\}.$$

Note that

$$\max\{e^{-2\kappa}, (t+1)^{-\frac{\kappa}{\kappa+1}}\} = \begin{cases} e^{-2\kappa}, & t \ge e^{2(\kappa+1)}-1 \\ (t+1)^{-\kappa}, & t \le e^{2(\kappa+1)}-1. \end{cases}$$

Hence for $\kappa < 1$, the number of steps required is $O\left(\frac{\ln(T+1)}{\kappa}\right)$ to get to a constant and $\frac{1}{\kappa}\ln\left(\frac{1}{t_1}\right)$ to get down to $t_1$ (as we need $s$ steps where $\left(1-\frac{\kappa}{\kappa+1}\right)^s \lesssim \ln(1+t_1) \sim t_1$), for a total of $O\left(\frac{1}{\kappa}\ln\left(\frac{T+1}{t_1}\right)\right)$. The last part follows from bounding every term and noting that $\sum_{k=1}^{N_{\text{step}}-1}\ln\left(\frac{t_{k+1}+1}{t_k+1}\right)\frac{1}{t_k+1}\varepsilon^2_{k+1} \le \sum_{k=1}^{N_{\text{step}}-1}\ln\left(\frac{t_{k+1}+1}{t_k+1}\right)\max_{1\le k\le N_{\text{step}}-1}\frac{1}{t_k+1}\varepsilon^2_{k+1} \lesssim \ln\left(\frac{T+1}{t_1+1}\right)\frac{\varepsilon^2}{\ln(T+1)} \le \varepsilon^2.$ □

Finally, we note that the TV distance between $p_0$ and $p_t$ for small $t$ can be bounded. This is true for any smooth $p_0$ [LLT23], but in our case admits a direct proof.

**Lemma 3.3.** *We have*
$$\mathrm{TV}(P_0, P_{\sigma^2}) \le \frac{1}{\sqrt{2}}\frac{\sigma^2\sqrt{n}}{\sigma^2_0}.$$

*Proof.* Note

$$\mathrm{KL}(\mathcal{N}(0,\sigma^2_0 I_n)\|\mathcal{N}(0,(\sigma^2_0+\sigma^2)I_n)) \le \frac{1}{2}\left[-n\ln\frac{\sigma^2_0+\sigma^2}{\sigma^2_0} - n + n\frac{\sigma^2_0+\sigma^2}{\sigma^2_0}\right]$$

$$\le \frac{n}{2}\left[-\frac{\sigma^2}{\sigma^2_0} + \frac{1}{2}\frac{\sigma^4}{\sigma^4_0} + \frac{\sigma^2}{\sigma^2_0}\right] \le \frac{n}{4}\frac{\sigma^4}{\sigma^4_0}.$$

By Pinsker's inequality,

$$\mathrm{TV}(\mathcal{N}(0,\sigma^2_0 I_n),\mathcal{N}(0,(\sigma^2_0+\sigma^2)I_n)) \le \sqrt{2\,\mathrm{KL}(\mathcal{N}(0,\sigma^2_0 I_n)\|\mathcal{N}(0,(\sigma^2_0+\sigma^2)I_n))} \le \frac{1}{\sqrt{2}}\frac{\sigma^2\sqrt{n}}{\sigma^2_0}.$$

The result follows from joint convexity of TV distance. □

## 3.2 Score function computation for Gaussian mixture

To learn a distribution using a standard diffusion model, we need to learn its score function (gradient of the log of the pdf) under Gaussian convolution, with varying levels of noise. For Gaussian mixtures, the score function has a particularly nice form, which we now derive.

First, we let $P_t = P_0 * \mathcal{N}(0, tI_n)$ and $Q_t = Q_0 * \mathcal{N}(0, tI_n)$. Then $P_t = Q_0 * \mathcal{N}(0, \sigma^2 I_n) = Q_{\sigma^2}$ where $\sigma^2 = \sigma_0^2 + t$. We consider the following probabilistic model where $\mu, \xi_1, \xi_2$ are drawn independently:

$$\mu \sim Q_0, \quad \xi_1, \xi_2 \sim \mathcal{N}(0, I_n), \quad X = \mu + \sigma_0 \xi_1, \quad Y = X + \sqrt{t}\xi_2 = \mu + \sigma_0 \xi_1 + \sqrt{t}\xi_2. \tag{5}$$

Then $X \sim P_0$ and $Y \sim P_t$. Letting $p_t, q_t$ be the corresponding densities and $V_t = \ln p_t$, we have by Tweedie's formula [Rob92] (see Appendix C for a derivation) that

$$\nabla V_t(y) = \frac{1}{\sigma^2}\mathbb{E}[\mu - y | Y = y] = \frac{1}{\sigma^2}\left(-y + \frac{\int_{\mathbb{R}^n} \mu \exp\left(\frac{\langle y, \mu\rangle}{\sigma^2} - \frac{|\mu|^2}{2\sigma^2}\right) dQ_0(\mu)}{\int_{\mathbb{R}^n} \exp\left(\frac{\langle y, \mu\rangle}{\sigma^2} - \frac{|\mu|^2}{2\sigma^2}\right) dQ_0(\mu)}\right) \tag{6}$$

$$f_{\sigma^2}(y) := y + \sigma^2 \nabla V_t(y) = \mathbb{E}[\mu | Y = y] = \frac{\int_{\mathbb{R}^n} \mu \exp\left(\frac{\langle y, \mu\rangle}{\sigma^2} - \frac{|\mu|^2}{2\sigma^2}\right) dQ_0(\mu)}{\int_{\mathbb{R}^n} \exp\left(\frac{\langle y, \mu\rangle}{\sigma^2} - \frac{|\mu|^2}{2\sigma^2}\right) dQ_0(\mu)}. \tag{7}$$

Note that in the same way we have

$$\nabla V_t(y) = \frac{1}{t}\mathbb{E}[X - y | Y = y] = \frac{1}{t}\left(-y + \frac{\int_{\mathbb{R}^n} x \exp\left(\frac{\langle y, x\rangle}{t} - \frac{|x|^2}{2t}\right) dP_0(x)}{\int_{\mathbb{R}^n} \exp\left(\frac{\langle y, x\rangle}{t} - \frac{|x|^2}{2t}\right) dP_0(x)}\right) \tag{8}$$

$$y + t\nabla V_t(y) = \mathbb{E}[X | Y = y] = \frac{\int_{\mathbb{R}^n} x \exp\left(\frac{\langle y, x\rangle}{t} - \frac{|x|^2}{2t}\right) dP_0(x)}{\int_{\mathbb{R}^n} \exp\left(\frac{\langle y, x\rangle}{t} - \frac{|x|^2}{2t}\right) dP_0(x)}. \tag{9}$$

Because of this identity, the score function can be learned as the minimal mean squared estimator in a (supervised) denoising problem of estimating $X$ given $Y$, known as the denoising auto-encoder (DAE) objective [Vin11]. (6)–(7) will be useful for analysis while (8)–(9) is used for the actual learning algorithm. Because of the nice form of (7), we will actually aim to learn $f_{\sigma^2}$. We remark that in the case of a finite mixture, (7) shows that $f_{\sigma^2}$ is represented by a softmax neural network with 1 hidden layer of $k$ units.

## 4 Learning the score for a single cluster

We will follow the approach in [KOS08], though with functions that are $\mathbb{R}^n$-valued rather than $\{0, 1\}$-valued.

First, we give some background on the Ornstein-Uhlenbeck process. Let $\gamma := \gamma_{\sigma^2}$ denote the density of $\mathcal{N}(0, \sigma^2 I_n)$. The generator $\mathscr{L} := \mathscr{L}_{\sigma^2}$ of the scaled Ornstein-Uhlenbeck process with variance $\sigma^2$, also known as Langevin dynamics for $\gamma$, is

$$\mathscr{L}f(x) = -\frac{1}{\sigma^2}\langle x, \nabla f(x)\rangle + \Delta f(x). \tag{10}$$

For a $\mathbb{R}^d$-valued function $f$, we interpret this componentwise, i.e.,

$$\mathscr{L}f(x) = -\frac{1}{\sigma^2}Df(x)x + \sum_{i=1}^{d}\partial_{ii}f.$$

The eigenfunctions are the (suitably scaled) Hermite polynomials $(h_{\mathbf{k}})_{\mathbf{k}\in\mathbb{N}_0^n}$, with $h_{\mathbf{k}} := h_{\mathbf{k},\sigma^2}$ having eigenvalue $-\frac{|\mathbf{k}|}{\sigma^2}$. Here we use $|\mathbf{k}|$ to denote $|\mathbf{k}|_1$. The Hermite polynomials form a complete eigenbasis for $L^2(\gamma)$. The scaled Ornstein-Uhlenbeck process can be described by the SDE

$$dx_t = -\frac{1}{\sigma^2}x_t + \sqrt{2}\,dW_t. \tag{11}$$

Let $(\mathscr{P}_t)_{t\geq 0}$ be the Markov semigroup of the scaled Ornstein-Uhlenbeck process, with generator given by (10). That is, we have

$$\mathscr{P}_t f(x) = \mathbb{E}_{x_0=x}f(x_t) \text{ when } x_t \text{ solves (11)}, \tag{12}$$
$$\mathscr{L}f = \lim_{t\to 0^+}\frac{\mathscr{P}_t f - f}{t}.$$

The following encapsulates the technique of noise sensitivity/stability: a bound on the $L^2(\gamma)$ norm of $\mathscr{L}^m f$ implies approximability of $f$ by a low-degree polynomial. To our knowledge, only the case $m = 1$ has been considered in the literature; however, similarly to how bounds on higher moments imply better tail bounds, bounding $\mathscr{L}^m f$ for larger $m$ can give better approximation.

**Lemma 4.1** (Noise stability implies low-degree approximability)**.** *Let $\gamma$ denote the density of $\mathcal{N}(0,\sigma^2 I_n)$. Suppose that $f : \mathbb{R}^d \to \mathbb{R}^d$ satisfies $\|f\|_{L^2(\gamma)} < \infty$, and that $m \in \mathbb{N}, L$ are such that $\|\mathscr{L}^m f\|_{L^2(\gamma)} \leq L^m$. For $d \in \mathbb{N}$, there exists a polynomial $g$ of degree $< d$ such that*

$$\|f - g\|_{L^2(\gamma)} \leq \left(\frac{L\sigma^2}{d}\right)^m.$$

*Proof.* Expand $f$ in the eigenfunction basis of $\mathscr{L}_{\sigma^2}$ as $f = \sum_{\mathbf{k}\in\mathbb{N}_0^d} a_{\mathbf{k}} h_{\mathbf{k}}$ where $a_{\mathbf{k}} \in \mathbb{R}^d$. Then

$$\mathscr{L}^d f = \sum_{\mathbf{k}\in\mathbb{N}_0^n} \frac{|\mathbf{k}|^d}{\sigma^{2d}} a_{\mathbf{k}} h_{\mathbf{k}}.$$

Hence

$$\frac{d^{2m}}{\sigma^{4m}}\sum_{|\mathbf{k}|\geq d}|a_{\mathbf{k}}|^2 \leq \sum_{\mathbf{k}\in\mathbb{N}_0^n}\frac{|\mathbf{k}|^{2m}}{\sigma^{4m}}|a_{\mathbf{k}}|^2 \leq L^{2m} \implies \sum_{|\mathbf{k}|\geq d}|a_{\mathbf{k}}|^2 \leq \left(\frac{L\sigma^2}{d}\right)^{2m}.$$

Taking $g = \sum_{|\mathbf{k}|<d} a_{\mathbf{k}} h_{\mathbf{k}}$, we have that $\|f - g\|_{L^2(\gamma)}^2 = \sum_{|\mathbf{k}|\geq d}|a_{\mathbf{k}}|^2$, which gives the desired bound. $\square$

## 4.1 Calculation of $\mathscr{L}^m f$

Let $\nu$ be a measure on $\mathbb{R}^n$ with all moments finite. Consider generating $y = x + \sigma\xi$ where $x \sim \nu$ and $\xi \sim \mathcal{N}(0, I_n)$. Let $\langle \cdot \rangle = \langle \cdot \rangle_y := \mathbb{E}_{P(x|y)}$ where $P(x|y)$ is the posterior distribution given by

$$\frac{dP(\cdot|y)}{d\nu}(x) \propto \exp\left(\frac{\langle y, x \rangle}{\sigma^2} - \frac{\|x\|^2}{2\sigma^2}\right).$$

(This is not to be confused with the inner product.) Then letting $f = f_{\sigma^2}$ as in (7), we have $f(y) := f_{\sigma^2} = \langle x \rangle_y$.

We first derive some formulas for differentiating posterior expectations $\langle \cdot \rangle$. Let $\tilde{x} = x - \langle x \rangle_y$ denote the centered random variable. For a function $g : \mathbb{R}^n \to \mathbb{R}$, we have the following general formula for differentiation with respect to $y_i$:

$$\sigma^2 \partial_i \langle g(x) \rangle_y = \frac{\int_{\mathbb{R}^n} g(x) x_i \exp\left(\frac{\langle y, x \rangle}{\sigma^2} - \frac{|x|^2}{2\sigma^2}\right) \nu(dx)}{\int_{\mathbb{R}^n} \exp\left(\frac{\langle y, x \rangle}{\sigma^2} - \frac{|x|^2}{2\sigma^2}\right) \nu(dx)}$$

$$- \frac{\int_{\mathbb{R}^n} g(x) \exp\left(\frac{\langle y, x \rangle}{\sigma^2} - \frac{|x|^2}{2\sigma^2}\right) \nu(dx) \int_{\mathbb{R}^n} x_i \exp\left(\frac{\langle y, x \rangle}{\sigma^2} - \frac{|x|^2}{2\sigma^2}\right) \nu(dx)}{\left(\int_{\mathbb{R}^n} \exp\left(\frac{\langle y, x \rangle}{\sigma^2} - \frac{|x|^2}{2\sigma^2}\right) \nu(dx)\right)^2}$$

$$= \langle g(x) \tilde{x}_i \rangle_y$$

$$\nabla \langle g(x) \rangle_y = \frac{1}{\sigma^2} \langle g(x) \tilde{x} \rangle_y.$$

More generally, let $\langle g(x^{(1)}, \ldots, x^{(r)}) \rangle_y$ denote $\mathbb{E}g(x^{(1)}, \ldots, x^{(r)})$ where $x^{(1)}, \ldots, x^{(r)}$ are independent draws from the posterior $p(\cdot|y)$. We use the trick of replacing the means with independent copies of the random variable (see e.g., [Tal10]) to obtain

$$\sigma^2 \nabla \left\langle g(x^{(1)}, \ldots, x^{(r)}) \right\rangle = \left\langle g(x^{(1)}, \ldots, x^{(r)})(\tilde{x}^{(1)} + \cdots + \tilde{x}^{(r)}) \right\rangle$$

$$= \left\langle g(x^{(1)}, \ldots, x^{(r)})(x^{(1)} + \cdots + x^{(r)} - rx^{(r+1)}) \right\rangle.$$

and for a function $h : \mathbb{R}^r \to \mathbb{R}$,

$$\nabla \left\langle g(x^{(1)}, \ldots, x^{(r)}) h\left(\left\langle x^{(1)}, y \right\rangle, \ldots, \left\langle x^{(r)}, y \right\rangle\right) \right\rangle$$

$$= \frac{1}{\sigma^2} \left\langle g(x^{(1)}, \ldots, x^{(r)}) h\left(\left\langle x^{(1)}, y \right\rangle, \ldots, \left\langle x^{(r)}, y \right\rangle\right)\left(x^{(1)} + \cdots + x^{(r)} - rx^{(r+1)}\right) \right\rangle$$

$$+ \left\langle g(x^{(1)}, \ldots, x^{(r)}) \sum_j \partial_j h\left(\left\langle x^{(1)}, y \right\rangle, \ldots, \left\langle x^{(r)}, y \right\rangle\right) x^{(j)} \right\rangle.$$

Now suppose $h$ is a homogeneous polynomial of degree $s$. Then

$$D \left\langle g(x^{(1)}, \ldots, x^{(r)}) h\left(\left\langle x^{(1)}, y \right\rangle, \ldots, \left\langle x^{(r)}, y \right\rangle\right) \right\rangle y$$

$$= \frac{1}{\sigma^2} \left\langle g(x^{(1)}, \ldots, x^{(r)}) h\left(\left\langle x^{(1)}, y \right\rangle, \ldots, \left\langle x^{(r)}, y \right\rangle\right)\left(\left\langle x^{(1)}, y \right\rangle + \cdots + \left\langle x^{(r)}, y \right\rangle - r\left\langle x^{(r+1)}, y \right\rangle\right) \right\rangle$$

$$+ s \left\langle g(x^{(1)}, \ldots, x^{(r)}) h\left(\left\langle x^{(1)}, y \right\rangle, \ldots, \left\langle x^{(r)}, y \right\rangle\right) \right\rangle \quad (13)$$

by Euler's formula $\sum_{j=1}^{r} \partial_j h(x_1, \ldots, x_j) x_j = s \cdot h$.

Now consider $h(z_1, \ldots, z_r) = \prod_{\ell=1}^{t} z_{j_\ell}$. Let $u(x^{(1)}, \ldots, x^{(r)}) = g(x^{(1)}, \ldots, x^{(r)}) h\left(\langle x^{(1)}, y \rangle, \ldots, \langle x^{(r)}, y \rangle\right)$. We have

$$
\Delta \left\langle u(x^{(1)}, \ldots, x^{(r)}) \right\rangle = \frac{1}{\sigma^4} \left\langle u \cdot \left\langle x^{(1)} + \cdots + x^{(r)} - r x^{(r+1)}, x^{(1)} + \cdots + x^{(r+1)} - (r+1) x^{(r+2)} \right\rangle \right\rangle
$$

$$
+ \frac{2}{\sigma^2} \left\langle g \cdot \sum_{\ell'=1}^{s} \prod_{\ell \neq \ell'} \left\langle x^{(i_\ell)}, y \right\rangle \cdot \left\langle x^{(i_{\ell'})}, x^{(1)} + \cdots + x^{(r)} - r x^{(r+1)} \right\rangle \right\rangle
$$

$$
+ \left\langle g \cdot \sum_{\substack{1 \leq \ell', \ell'' \leq s \\ \ell' \neq \ell''}} \prod_{\ell \neq \ell', \ell''} \left\langle x^{(i_\ell)}, y \right\rangle \cdot \left\langle x^{(i_{\ell'})}, x^{(i_{\ell''})} \right\rangle \right\rangle \quad (14)
$$

Applying the above to monomials $g(x^{(1)}, \ldots, x^{(r)}) = \prod_{\ell=1}^{s} \left\langle x^{(i_\ell)}, x^{(i'_\ell)} \right\rangle$ and using induction, we have the following.

**Lemma 4.2.** *Let $f = f_{\sigma^2}$ be as in (7). We have*

$$
\mathscr{L}^m f(y) = \left\langle x^{(1)} \sum_{s+t \leq m} \sum_{i, i' \in [2m+1]^s, j \in [2m+1]^t} a_{i, i', j} \sigma^{-2(s+t+m)} \prod_{\ell=1}^{s} \left\langle x^{(i_\ell)}, x^{(i'_\ell)} \right\rangle \prod_{\ell=1}^{t} \left\langle x^{(j_\ell)}, y \right\rangle \right\rangle \quad (15)
$$

*where $\sum_{i, i', j} |a_{i, i', j}| \leq 30^m m!^2$.*

*Proof.* Recall that $\mathscr{L} f(x) = -\frac{1}{\sigma^2} Df(x)x + \Delta f(x)$. We induct on $m$. Suppose that the lemma holds for $m - 1$; each term has at most a number of replicas $r \leq 2m - 1$ and $s, t \leq m - 1$. We consider the effect of the map $f \mapsto \frac{1}{\sigma^2} Df(y) \cdot y$ (13) and $\Delta$ (14) on a single term in (15). First, note that in each case, we obtain a factor $\frac{1}{\sigma^2}$, and a factor of $\frac{1}{\sigma^2}$ for each additional $\left\langle x^{(i_\ell)}, x^{(i'_\ell)} \right\rangle$ term as well as $\left\langle x^{(i_\ell)}, y \right\rangle$ term (with a $\sigma^2$ factor if we remove a term); this justifies the $\sigma^{-2(s+t+m)}$ factor.

1. In (13), the number of replicas $r$ in the resulting terms increases by at most 1, $t$ increases by at most 1, and the sum of absolute value of coefficients is at most

$$
2r + s \leq 2(2m - 1) + (m - 1) \leq 5m - 3.
$$

2. In (14), the number of replicas $r$ in the resulting terms increases by at most 2, $s$ increases by at most 1, and the sum of absolute value of coefficients is at most

$$
2r \cdot 2(r + 1) + 2 \cdot s \cdot 2r + s(s - 1)
$$
$$
\leq 4 \cdot (2m - 1)(2m + 1) + 2 \cdot (m - 1) \cdot 2(2m - 1) + (m - 1)(m - 2)
$$
$$
\leq 16m^2 + 8m^2 + m^2 \leq 25m^2
$$

The sum of absolute value of coefficients multiplies by at most $30m^2$. This finishes the induction step. $\square$

## 4.2 Bounding $\mathscr{L}^m f$

We would like to bound $\|\mathscr{L}^m f\|_{L^p(\gamma_{\sigma^2})}$ for all $m$. The kind of growth we get in $m$ is captured by the following definition.

**Definition 4.3.** *We say $f : \mathbb{R}^n \to \mathbb{R}^{n'}$ is a $(r, \sigma)$-**Gaussian-noise-sensitive function** if for all $m \in \mathbb{N}$ and $p \geq 1$,*

$$\|\mathscr{L}^m f\|_{L^p(\gamma_{\sigma^2})} \leq r\sigma \left(\frac{rm^2}{\sigma^2}\right)^m \max\{r, \sqrt{mp}\}^m.$$

**Lemma 4.4.** *Let $f = f_{\sigma^2}$ be as in (7). Suppose $Q_0$ is supported on $B_R(0)$, $R \geq \sigma$. Let $P$ be the density of $Q_0 * \mathcal{N}(0, \sigma^2 I_n)$. Then for any $m \in \mathbb{N}$, we have the following:*

$$\|\mathscr{L}^m f\|_{L^p(P)} = R \cdot O\left(\left(\frac{1}{\sigma}\right)^2 \left(\frac{mR}{\sigma}\right)^2 \left(1 + \frac{(mp)^{1/2}\sigma}{R}\right)\right)^m$$

$$\|\mathscr{L}^m f\|_{L^p(\gamma_{\sigma^2})} \leq R \cdot O\left(\left(\frac{1}{\sigma}\right)^2 \left(\frac{m^2 R}{\sigma}\right) \max\left\{\frac{R}{\sigma}, (mp)^{1/2}\right\}\right)^m$$

*Therefore, $f$ is $\left(O\left(\frac{R}{\sigma}\right), \sigma\right)$-Gaussian-noise-sensitive.*

*Proof.* In (15) in Lemma 4.2,

$$\left|\prod_{\ell=1}^{s} \left\langle x^{(i_\ell)}, x^{(i'_\ell)} \right\rangle\right| \leq R^{2s} \leq R^{2m}.$$

Note that the joint distribution of $(x^{(j)}, y)$ is the same for any $j$, namely, it is the distribution when $x^{(j)} \sim p_0$ and $y = x^{(j)} + \sigma\xi^{(j)}$ when $\xi \sim \mathcal{N}(0, I_n)$ is independent of $x^{(j)}$. Then for $p \in \mathbb{N}$,

$$\mathbb{E}\left\langle \left|\prod_{\ell=1}^{t} \left\langle x^{(j_\ell)}, y \right\rangle\right|\right\rangle^p$$

$$\leq \mathbb{E}_{\substack{x^{(1)} \sim p_0,\, y = x^{(1)} + \sigma\xi \\ x^{(r)} \sim p(\cdot|y),\, r > 1}} \left|\prod_{\ell=1}^{t} \left\langle x^{(j_\ell)}, y \right\rangle\right|^p \qquad\qquad \text{by Jensen's inequality}$$

$$\leq \prod_{\ell=1}^{t} \left[\mathbb{E}\left|\left\langle x^{(j_\ell)}, y \right\rangle\right|^{tp}\right]^{1/t} \qquad\qquad \text{by Hölder's inequality}$$

$$\leq \mathbb{E}\left|\left\langle x^{(1)}, x^{(1)} + \sigma\xi^{(1)} \right\rangle\right|^{tp} \qquad\qquad (x^{(j)}, y) \overset{d}{=} (x^{(1)}, x^{(1)} + \sigma\xi^{(1)})$$

$$\leq \mathbb{E}\sum_{k=0}^{tp} \binom{tp}{k} \left|x^{(1)}\right|^{2(tp-k)} \sigma^k \left|\left\langle x^{(1)}, \xi^{(1)} \right\rangle\right|^k \qquad\qquad \text{Binomial theorem}$$

$$\leq \sum_{k=0}^{tp} \binom{tp}{k} R^{2(tp-k)} \sigma^k R^k \mathbb{E}_{X \sim \mathcal{N}(0,1)}|X|^k$$

$$\leq \sum_{k=0}^{tp} \binom{tp}{k} R^{2(tp-k)} \sigma^k R^k (k-1)!! \qquad\qquad \text{Gaussian moment bound}$$

$$\leq \sum_{k=0}^{tp} \binom{tp}{k} R^{2(tp-k)} \sigma^k R^k (tp)^{k/2} \leq R^{2tp} \left( 1 + \frac{(tp)^{1/2}\sigma}{R} \right)^{tp}.$$

Consider one term $a_{i,i',j} \left\langle x^{(1)} \sigma^{-2(s+t+m)} \prod_{\ell=1}^s \left\langle x^{(i_\ell)}, x^{(i'_\ell)} \right\rangle \prod_{\ell=1}^t \left\langle x^{(j_\ell)}, y \right\rangle \right\rangle$ in (15). We have by Jensen's inequality that

$$\left\| \left\langle x^{(1)} \sigma^{-2(s+t+m)} \prod_{\ell=1}^s \left\langle x^{(i_\ell)}, x^{(i'_\ell)} \right\rangle \prod_{\ell=1}^t \left\langle x^{(j_\ell)}, y \right\rangle \right\rangle \right\|_{L^p(P)}$$

$$\leq \left[ \mathbb{E} \left\langle \left| x^{(1)} \sigma^{-2(s+t+m)} \prod_{\ell=1}^s \left\langle x^{(i_\ell)}, x^{(i'_\ell)} \right\rangle \prod_{\ell=1}^t \left\langle x^{(j_\ell)}, y \right\rangle \right|^p \right\rangle \right]^{\frac{1}{p}}$$

$$\leq \sigma^{-2(s+t+m)} R^{2s+1} R^{2t} \left( 1 + \frac{(tp)^{1/2}\sigma}{R} \right)^t$$

$$= O \left( \frac{R}{\sigma^{2m}} \left( \frac{R}{\sigma} \right)^{2(s+t)} \left( 1 + \frac{(tp)^{1/2}\sigma}{R} \right)^t \right) = O \left( \frac{R}{\sigma^{2m}} \left( \frac{R}{\sigma} \right)^{2m} \left( 1 + \frac{(mp)^{1/2}\sigma}{R} \right)^m \right)$$

using $s + t \leq m$. Then

$$\|\mathscr{L}^m f\|_{L^p(P)} \leq 30^m m!^2 \cdot O \left( \frac{R}{\sigma^{2m}} \left( \frac{R}{\sigma} \right)^{2m} \left( 1 + \frac{(mp)^{1/2}\sigma}{R} \right)^m \right)$$

$$= R \cdot O \left( \left( \frac{1}{\sigma} \right)^{2m} \left( \frac{mR}{\sigma} \right)^{2m} \left( 1 + \frac{(mp)^{1/2}\sigma}{R} \right)^m \right) \tag{16}$$

By Hölder's inequality and Lemma A.4 (appropriately scaled),

$$\|\mathscr{L}^m f\|_{L^p(\gamma_{\sigma^2})} = \left( \int_{\mathbb{R}^n} |\mathscr{L}^m f|^p \frac{d\gamma_{\sigma^2}}{dP} dP \right)^{\frac{1}{p}} \leq \|\mathscr{L}^m f\|_{L^{p(1+q)}(\gamma_{\sigma^2})} \left\| \frac{d\gamma_{\sigma^2}}{dP} \right\|_{L^{1+\frac{1}{q}}(P)}^{\frac{1}{p}}$$

$$\leq O \left( \left( \frac{1}{\sigma} \right)^2 \left( \frac{mR}{\sigma} \right)^2 \left( 1 + \frac{(mp(q+1))^{1/2}\sigma}{R} \right) \right)^m e^{\frac{(R/\sigma)^2}{2pq}}.$$

To optimize this bound, set $q = \frac{(R/\sigma)^2}{pm}$. Then by Hölder's inequality,

$$\|\mathscr{L}^m f\|_{L^p(\gamma_{\sigma^2})} \leq Re^{\frac{(R/\sigma)^2}{2pq}} \cdot O \left( \left( \frac{1}{\sigma} \right)^2 \left( \frac{mR}{\sigma} \right)^2 \left( 1 + \frac{mp(q+1)\sigma}{R} \right) \right)^m$$

$$\leq Re^{m/2} \cdot O \left( \left( \frac{1}{\sigma} \right)^2 \left( \frac{mR}{\sigma} \right)^2 \right)^m \cdot \left[ 1 + O \left( \frac{(mpq)^{1/2}\sigma}{R} \right)^m + O \left( \frac{(mp)^{1/2}\sigma}{R} \right)^m \right]$$

The first two terms give $R \cdot O \left( \left( \frac{1}{\sigma} \right)^2 \left( \frac{mR}{\sigma} \right)^2 \right)^m$, while the last term gives $R \cdot O \left( \left( \frac{1}{\sigma} \right)^2 \left( \frac{m^2 R}{\sigma} \right) (mp)^{1/2} \right)$. This completes the proof. $\qquad \square$

## 4.3 Approximation with polynomial

Suppose we want to approximate $f = f_{\sigma^2}$ (in (7)) with a low-degree polynomial on the mixture distribution $\gamma'$. We seek a low-degree polynomial $g$ such that $\|f - g\|^2_{\gamma'}$ is small. To do this, we first smooth $f$ by the Ornstein-Uhlenbeck semigroup to obtain $\mathscr{P}_t f$ (see (12)), and then find a $g$ that approximates $\mathscr{P}_t f$. We can bound using Cauchy-Schwarz that

$$\|f - g\|^2_{\gamma'} \leq 2\left(\|f - \mathscr{P}_t f\|^2_{\gamma'} + \|\mathscr{P}_t f - g\|^2_{\gamma'}\right)$$
$$\leq 2\left(\|f - \mathscr{P}_t f\|^2_{\gamma'} + \|\mathscr{P}_t f - g\|^2_{L^2(\gamma)} + \|\mathscr{P}_t f - g\|^2_{L^4(\gamma)}\, \chi^2(\gamma'\|\gamma)^{1/2}\right).$$

Doing the smoothing ensures that we have better control over the term with higher $p$-norm, $\|\mathscr{P}_t f - g\|_{L^4(\gamma)}$. Choosing $t = \varepsilon$, $\mathscr{P}_t f$ has exponential decay in coefficients with rate $\varepsilon$, so we can approximate $\mathscr{P}_t f$ with a degree-$\Theta\left(\frac{1}{\varepsilon}\right)$ polynomial. To bound $\|f - \mathscr{P}_t f\|^2_{\gamma'}$, we bound its derivative:

$$\frac{d}{dt}\|f - \mathscr{P}_t f\|^2_{\gamma'} \leq \int_{\mathbb{R}^n} 2(f - \mathscr{P}_t f)(-\mathscr{L}\mathscr{P}_t f)\, d\gamma'$$
$$\leq 2\|f - \mathscr{P}_t f\|_{\gamma'}\,\|\mathscr{L}\mathscr{P}_t f\|_{\gamma'}$$
$$\implies \frac{d}{dt}\|f - \mathscr{P}_t f\|_{\gamma'} \leq \|\mathscr{L}\mathscr{P}_t f\|_{\gamma'}.$$

Focusing on the dependence of total error on $\varepsilon$, we can bound this with a change-of-measure inequality and Lemma 4.4, which, after integrating in $t$, gives error $O(\varepsilon)$. In order to obtain poly-logarithmic degree, we need a higher-order version of this argument. In preparation for Section 5, when we need to consider the norm with respect to a different measure, we state the following more generally.

**Lemma 4.5.** *Suppose $f$ is $(O(r), \sigma)$-Gaussian-noise-stable and $\gamma'$ is a measure such that for all $a \geq 0$, $\left\|\frac{d\gamma'}{d\gamma}\right\|_{L^{1+a}(\gamma)} \leq e^{\frac{ar^2}{2}}$ (for $\gamma := \gamma_{\sigma^2}$) (e.g., from Lemma A.4, $Q * \mathcal{N}(0, I_n)$ where $Q$ is supported on $B_r(0)$), where $r \geq 1$. There is a polynomial $g$ of degree at most $O\left(r^4 \ln\left(\frac{1}{\varepsilon}\right)^4 \max\left\{r^2, \ln\left(\frac{1}{\varepsilon}\right)\right\}\right)$ such that*

$$\|f - g\|^2_{\gamma'} \leq \varepsilon^2 r^2 \sigma^2 \quad and \quad \|g\|^2_\gamma \leq \|f\|^2_\gamma.$$

*Proof.* To get a better bound, we use a more clever smoothing strategy. For some $\widetilde{f}$ to be defined, we will bound

$$\|f - g\|^2_{\gamma'} \leq 2\left(\left\|f - \widetilde{f}\right\|^2_{\gamma'} + \left\|\widetilde{f} - g\right\|^2_{\gamma'}\right), \tag{17}$$

where $g$ is a polynomial approximation of $\widetilde{f}$ obtained by truncating the Hermite expansion.

To define $\widetilde{f}$, we approximate $f$ with a numerical differentiation formula for a higher derivative. Define the finite difference by $\Delta_{x,h}g(x) = \Delta_h g(x) := g(x + h) - g(x)$. We can write this as $\Delta_h g = T_h g - g$, where $T_h g(x) := g(x + h)$. Suppose that $g \in C^m$. By the Binomial Theorem on

$T_h$ – id and Taylor's Theorem,

$$\Delta_h^m g(0) = \sum_{j=0}^m \binom{m}{j} (-1)^{m-j} \left[ \sum_{i=0}^{m-1} g^{(i)}(0) \frac{(hj)^i}{i!} + \frac{g^{(m)}(\xi_j)}{m!} (hj)^m \right] \text{ for some } \xi_j \in [0, hj]$$

$$= \sum_{i=0}^{m-1} \frac{g^{(i)}(0)}{i!} \Delta_{x,h}^m (x^i)|_{x=0} + h^m \sum_{j=0}^m \binom{m}{j} (-1)^{m-j} \frac{g^{(j)}(\xi_j)}{m!} j^m$$

$$= h^m \sum_{j=0}^m \binom{m}{j} (-1)^{m-j} \frac{g^{(j)}(\xi_j)}{k!} j^m$$

where in the last step we use the fact that finite differencing reduces the degree of a polynomial by 1, the $m$th finite difference of polynomials of degree $< m$ is 0. Hence

$$|\Delta_h^m g(0)| \le h^m \sum_{j=0}^m \binom{m}{j} \left( \frac{ej}{m} \right)^m \max_{\xi \in [0,kh]} |f^{(m)}(\xi)| \le h^m (2e)^m \max_{\xi \in [0,kh]} |f^{(m)}(\xi)|.$$

Let $\widetilde{f} = \sum_{j=1}^m (-1)^{j+1} \binom{m}{j} P_{jh} f$. Note $f - \widetilde{f} = (-1)^m \Delta_{t,h}^m (\mathscr{P}_t f)|_{t=0}$. We then have (noting that $\mathscr{P}_s$ and $\mathscr{L}$ commute)

$$\left\| f - \widetilde{f} \right\|_{\gamma'}^2 \le (2he)^{2m} \int_{\mathbb{R}^n} \max_{t \in [0,mh]} \left( \frac{d^m}{dt^m} \mathscr{P}_t f(x) \right)^2 d\gamma'(x)$$

$$\le (2he)^{2m} \int_{\mathbb{R}^n} \left( \frac{d^m}{dt^m} \mathscr{P}_t f(x)|_{t=0} + \int_0^{mh} \left| \frac{d^{m+1}}{ds^{m+1}} \mathscr{P}_s f(x) \right| ds \right)^2 d\gamma'(x)$$

$$\le 2(2he)^{2m} \left( \int_{\mathbb{R}^n} |\mathscr{L}^m f(x)|^2 d\gamma'(x) + mh \int_0^{mh} \int_{\mathbb{R}^n} |\mathscr{P}_s \mathscr{L}^{m+1} f(x)|^2 d\gamma'(x) \, ds \right). \quad (18)$$

First we bound the first term and the second term when $s = 0$. We use the Gaussian noise sensitivity assumption, the assumption on $\gamma'$, and Hölder's inequality to obtain

$$\int_{\mathbb{R}^n} |\mathscr{L}^l f|^2 \, d\gamma' \le \left\| \mathscr{L}^l f \right\|_{L^{(p+1)}(\gamma)}^2 \left\| \frac{d\gamma'}{d\gamma} \right\|_{L^{\frac{p+1}{p}}(\gamma)}$$

$$\le r^2 \sigma^2 O \left( \frac{rl^2}{\sigma^2} \right)^{2l} \max\{r^2, l(1+p)\}^l e^{\frac{r^2}{2p}}. \quad (19)$$

To get a bound for the second term when $s > 0$, we bound the following derivative, again using the Gaussian noise sensitivity assumption and Hölder's inequality:

$$\frac{d}{ds} \int_{\mathbb{R}^n} |\mathscr{P}_s \mathscr{L}^l f|^2 \, d\gamma' \le \int_{\mathbb{R}^n} 2 \left\langle \mathscr{P}_s \mathscr{L}^l f, \mathscr{P}_s \mathscr{L}^{l+1} f \right\rangle \frac{d\gamma'}{d\gamma} \, d\gamma$$

$$\le 2 \left\| \mathscr{P}_s \mathscr{L}^l f \right\|_{L^{2(p+1)}(\gamma)} \left\| \mathscr{P}_s \mathscr{L}^{l+1} f \right\|_{L^{2(p+1)}(\gamma)} \left\| \frac{d\gamma'}{d\gamma} \right\|_{L^{\frac{p+1}{p}}(\gamma)}$$

$$\le 2 \left\| \mathscr{L}^l f \right\|_{L^{2(p+1)}(\gamma)} \left\| \mathscr{L}^{l+1} f \right\|_{L^{2(p+1)}(\gamma)} \left\| \frac{d\gamma'}{d\gamma} \right\|_{L^{\frac{p+1}{p}}(\gamma)} \quad (20)$$

$$\le r^2 \sigma^2 O \left( \frac{rl^2}{\sigma^2} \right)^{2l} \max\{r^2, l(1+p)\}^l e^{\frac{r^2}{2p}}, \quad (21)$$

23

where in (20) we use the fact that for $q \geq 1$, $\|\mathscr{P}_s g\|_{L^q(\gamma)}$ is monotonically decreasing in $s$. In both (19) and (21), we optimize the bound by taking $p = \frac{r^2}{l}$ to obtain

$$\int_{\mathbb{R}^n} |\mathscr{L}^l f|^2 \, d\gamma' \leq r^2 \sigma^2 O\left(\frac{rl^2}{\sigma^2}\right)^{2l} \max\{r^2, l\}^l, \qquad\qquad l = m, \, m+1$$

$$\frac{d}{ds} \int_{\mathbb{R}^n} |\mathscr{L}^l \mathscr{P}_s f|^2 \, d\gamma' \leq r^2 \sigma^2 O\left(\frac{rl^2}{\sigma^2}\right)^{2l} \max\{r^2, l\}^l, \qquad\qquad l = m+1.$$

Integrating the last inequality twice and substituting into (18) gives

$$\left\|f - \widetilde{f}\right\|_{\gamma'}^2 \lesssim r^2 \sigma^2 O(h)^{2m} \left((1 + (mh)^3)\left(\frac{rm^2}{\sigma^2}\right)^{2m} \max\{r^2, m\}^m\right)$$

$$\lesssim r^2 \sigma^2 O\left(\frac{hrm^2 \max\{r, \sqrt{m}\}}{\sigma^2}\right)^{2m} \tag{22}$$

when $mh = O(1)$. We choose $m \sim \ln\left(\frac{1}{\varepsilon}\right)$ and $h \leq \frac{c\sigma^2}{rm^2 \max\{r, \sqrt{m}\}}$ for an appropriate constant $c$ to get

$$\left\|f - \widetilde{f}\right\|_{\gamma'}^2 \lesssim \varepsilon^2 r^2 \sigma^2.$$

Write $f = \sum_{\mathbf{k} \in \mathbb{N}_0^d} a_{\mathbf{k}} h_{\mathbf{k}}$ where $a_{\mathbf{k}} \in \mathbb{R}^d$, and let $p_l = \sum_{|\mathbf{k}|=l} a_{\mathbf{k}} h_{\mathbf{k}}$, so that $f = \sum_{l=0}^{\infty} p_l$. Then

$$\widetilde{f} = \sum_{\mathbf{k} \in \mathbb{N}_0^n} b_{\mathbf{k}} h_{\mathbf{k}} \quad \text{for} \quad b_{\mathbf{k}} = \sum_{j=1}^{m} (-1)^{j+1} \binom{m}{j} e^{-jh|\mathbf{k}|/\sigma^2} a_{\mathbf{k}} = [1 - (1 - e^{-h|\mathbf{k}|/\sigma^2})^m] a_{\mathbf{k}}. \tag{23}$$

Let $q_l = \sum_{|\mathbf{k}|=l} b_{\mathbf{k}} h_{\mathbf{k}}$. Note

$$\|q_l\|_\gamma^2 = \sum_{|\mathbf{k}|=l} |b_{\mathbf{k}}|^2 = \sum_{|\mathbf{k}|=l} [1 - (1 - e^{-h|\mathbf{k}|/\sigma^2})^m]^2 |a_{\mathbf{k}}|^2 \leq m^2 e^{-2hl/\sigma^2} \|p_l\|_\gamma^2.$$

We approximate $\widetilde{f}$ with $g = \sum_{|\mathbf{k}|<L} b_{\mathbf{k}} h_{\mathbf{k}}$ where $L \geq Ch^{-1} \ln\left(\frac{1}{\varepsilon}\right)$ for an appropriate constant $C$. We have by Cauchy-Schwarz and Lemma A.3 that

$$\left\|\widetilde{f} - g\right\|_{\gamma'}^2 \leq \int_{\mathbb{R}^n} \left|\sum_{l \geq L} q_l\right|^2 \, d\gamma'$$

$$\leq \sum_{l \geq L} e^{-hl} \cdot \sum_{l \geq L} e^{hl} \int_{\mathbb{R}^n} |q_l|^2 \, d\gamma'$$

$$\leq 1 \cdot \sum_{l \geq L} e^{hl} \|q_l\|_\gamma^2 e^{2\sqrt{l}r}$$

$$\leq \sum_{l \geq L} e^{hl} m^2 e^{-2hl/\sigma^2} \|p_l\|_\gamma^2 e^{2\sqrt{l}r}$$

$$\leq \max_{l \geq L} m^2 e^{2\sqrt{l}r - hl/\sigma^2} \|f\|_\gamma^2 \leq \varepsilon^2 \|f\|_\gamma^2 \tag{24}$$

24

when we take $L \geq C\frac{\sigma^4 r^2}{h^2} \asymp r^4 \ln\left(\frac{1}{\varepsilon}\right)^4 \max\left\{r^2, \ln\left(\frac{1}{\varepsilon}\right)\right\}$ for an appropriate constant $C$. Note that the Gaussian-noise-sensitive property also implies $\|f\|_\gamma^2 \leq r^2\sigma^2$. Plugging into (17) the inequalities (22) and (24),

$$\|f - g\|_{\gamma'}^2 \lesssim \left\|f - \tilde{f}\right\|_{\gamma'}^2 + \left\|\tilde{f} - g\right\|_{\gamma'}^2 \lesssim r^2\sigma^2\varepsilon^2.$$

Choosing constants appropriately then gives the desired bound. Finally note that Equation (23) implies $|a_{\mathbf{k}}| \leq |b_{\mathbf{k}}|$ for all $\mathbf{k} \in \mathbb{N}_0^n$. Hence

$$\|g\|_\gamma^2 \leq \left\|\tilde{f}\right\|_\gamma^2 = \sum_{\mathbf{k} \in \mathbb{N}_0^n} |b_{\mathbf{k}}|^2 \leq \sum_{\mathbf{k} \in \mathbb{N}_0^n} |a_{\mathbf{k}}|^2 = \|f\|_\gamma^2.$$

$\square$

# 5 From one to multiple clusters

**Definition 5.1:** Let $\mathcal{C} \subset \mathbb{R}^n$. We say $\mathcal{C}$ is a **complete set of $R$-warm starts** for $Q_0$ if

$$Q_0\left(\bigcup_{\widehat{\mu} \in \mathcal{C}} B_R(\widehat{\mu})\right) = 1.$$

In other words, for all $\mu$ in the support of $Q_0$, there exists $\widehat{\mu} \in \mathcal{C}$ such that $\|\widehat{\mu} - \mu\| \leq R$.

**Definition 5.2:** Let $\mathcal{C} = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_{k'}\}$. The **Voronoi partition** $V_1, \ldots, V_{k'}$ corresponding to $\mathcal{C}$ is defined by $V_j = \left\{x \in \mathbb{R}^n : |x - \widehat{\mu}_j| = \min_{1 \leq j' \leq k'} |x - \widehat{\mu}_{j'}|\right\}$.

Note that up to the boundaries (which are a measure 0 set), this induces a partition of $\mathbb{R}^n$.

To tackle the multiple cluster setting, we first suppose that we have a complete set of $R$-warm starts $\mathcal{C} = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_{k'}\}$, and let $V_1, \ldots, V_{k'}$ be the corresponding Voronoi partition. For a probability measure $P$ and set $S$, define the unnormalized restriction by $P_S(A) = P(A \cap S)$ and the normalized restriction by $P|_S(A) = \frac{P(A \cap S)}{P(S)}$. Define $P^S = (Q_0)_S * \mathcal{N}(0, \sigma^2 I_n)$—that is, we do the restriction before the convolution. Define

$$f_{S,\sigma^2}(y) = \mathbb{E}_{\substack{\mu \sim Q_0|_S \\ Y = \mu + \sigma\xi,\, \xi \sim \mathcal{N}(0, I_n)}} [\mu \mid Y = y] = y + \sigma^2 \nabla \ln p^S(y).$$

When $\sigma$ is understood, we omit it from the subscript.

We need to show that we still have a good polynomial approximation for $f$ under the measure $P|_{V_i} = (Q_0 * \mathcal{N}(0, \sigma^2 I_n))|_{V_i}$. By the analysis for one cluster, we have good approximation of $f_{V_i,\sigma^2}$ under $P^{V_i}$. To obtain the result for multiple clusters, we need to show that this approximation is preserved even when we consider $f$ instead of $f_{V_i,\sigma^2}$ and $P|_{V_i}$ instead of $P^{V_i}$, i.e., deal with the leakage into $V_i$ from the other Voronoi cells after $Q_0$ is convolved with $\mathcal{N}(0, \sigma^2 I)$, in both the score function and the measure which the norm is with respect to. We do this in Sections 5.1 and 5.2, respectively.

## 5.1 Estimation with the "nearby" score

We will actually bound $\|f - f_{S_i}\|_{L^2(P_{V_i})}$, for an expanded neighborhood $S_i$ of $\widehat{\mu}_i$, allowing an extra "buffer region" where mass is allowed to leak in.

**Lemma 5.3.** *Suppose that $\mathcal{C} = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_{k'}\}$ is a complete set of $R$-warm starts for $Q_0$, and let $V_1, \ldots, V_{k'}$ be the corresponding Voronoi partition. Suppose $Q_0$ is supported on $B_M(0)$. Given $R' > R$, let $S_i = B_{R'}(\widehat{\mu}_i)$ and define*

$$f^{\mathcal{C},R'}_{\mathrm{loc},\sigma^2}(y) = f_{S_i,\sigma^2}(y) = \mathbb{E}_{\substack{\mu \sim Q_0|_{S_i} \\ Y = \mu + \sigma\xi, \, \xi \sim \mathcal{N}(0, I_n)}} [\mu \mid Y = y] = y + \sigma^2 \nabla \ln p^{S_i}(y) \text{ when } y \in V_i.$$

*We write $f_{\mathrm{loc},\sigma^2}$ when $\mathcal{C}, R'$ are clear. Then for $R' = 3R + 2\sqrt{2}\sigma\sqrt{\ln\left(\frac{k'}{\varepsilon}\right)}$,*

$$\left\|f_{\sigma^2} - f_{\mathrm{loc},\sigma^2}\right\|^2_{L^2(P)} \le \left(32\sigma^2 + 6R'^2\right)\varepsilon \lesssim \left(R^2 + \sigma^2 \ln\left(\frac{k'}{\varepsilon}\right)\right)\varepsilon.$$

Note that

$$\left\|f_{\sigma^2} - f_{\mathrm{loc},\sigma^2}\right\|^2_{L^2(P)} = \sum_{i=1}^{k'} \left\|f_{\sigma^2} - f_{S_i,\sigma^2}\right\|^2_{L^2(P_{V_i})},$$

so this gives a bound for the $L^2$ norm within each Voronoi cell.

To prove this, we first show that with high probability, $Y = \mu + \sigma\xi$ does not stray too far from the Voronoi cell of $\mu$.

**Lemma 5.4.** *Let $\mathcal{C} = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_{k'}\}$ and $V_1, \ldots, V_{k'}$ be the corresponding Voronoi partition. Suppose that $|\mu - \widehat{\mu}_i| \le R$. Define $Y = \mu + \sigma\xi$ where $\xi \sim \mathcal{N}(0, I_n)$, and let $i'$ be such that $Y \in V_{i'}$. Then with probability $\ge 1 - \varepsilon$,*

$$|\mu - \widehat{\mu}_{i'}| \le 3R + 2\sqrt{2}\sigma\sqrt{\ln\left(\frac{k'}{\varepsilon}\right)}.$$

In other words, with high probability, even if $Y = \mu + \sigma\xi \in V_{i'}$ for $i' \ne i$ (i.e., adding a Gaussian brings the point to a different Voronoi cell), $\mu$ will not be too far from the center of the new cell $\widehat{\mu}_{i'}$.

*Proof.* Let $v_{ii'} = \frac{\widehat{\mu}_{i'} - \widehat{\mu}_i}{|\widehat{\mu}_{i'} - \widehat{\mu}_i|}$ be the unit vector pointing from $\widehat{\mu}_i$ toward $\widehat{\mu}_{i'}$. Since $\mu \in V_i$, we have

$$\langle \mu - \widehat{\mu}_i, v_{ii'} \rangle \le \frac{|\widehat{\mu}_{i'} - \widehat{\mu}_i|}{2}.$$

By Gaussian tail bounds and a union bound, with probability $\ge 1 - \varepsilon$, we have that for all $i'$,

$$\langle Y - \widehat{\mu}_i, v_{ii'} \rangle = \langle \mu + \sigma\xi - \widehat{\mu}_i, v_{ii'} \rangle \le \langle \mu - \widehat{\mu}_i, v_{ii'} \rangle + \sqrt{2}\sigma\sqrt{\ln\left(\frac{k'}{\varepsilon}\right)} \le R + \sqrt{2}\sigma\sqrt{\ln\left(\frac{k'}{\varepsilon}\right)}$$

Therefore, if $y \in V_{i'}$, then we have

$$\frac{|\widehat{\mu}_{i'} - \widehat{\mu}_i|}{2} \le \langle Y - \widehat{\mu}_i, v_{ii'} \rangle \le R + \sqrt{2}\sigma\sqrt{\ln\left(\frac{k'}{\varepsilon}\right)}$$

and

$$|\mu - \widehat{\mu}_{i'}| \le |\mu - \widehat{\mu}_i| + |\widehat{\mu}_i - \widehat{\mu}_{i'}| \le 3R + 2\sqrt{2}\sigma\sqrt{\ln\left(\frac{k'}{\varepsilon}\right)}.$$

$\square$

*Proof of Lemma 5.3.* Given $\mu$, let $i$ be such that $y \in V_i$. We note

$$f_{\sigma^2}(y) = \mathbb{E}[\mathbb{1}_{S_i}(\mu)\mu \mid Y = y] + \mathbb{E}[\mathbb{1}_{S_i^c}(\mu)\mu \mid Y = y]$$

$$f_{\mathrm{loc},\sigma^2}(y) = \mathbb{P}(\mu \in S_i \mid Y = y)f_{\mathrm{loc},\sigma^2}(y) + \mathbb{P}(\mu \notin S_i \mid Y = y)f_{\mathrm{loc},\sigma^2}(y)$$

$$= \mathbb{P}(\mu \in S_i \mid Y = y)\frac{\mathbb{E}\left[\mathbb{1}_{S_i}(\mu)\mu \mid Y = y\right]}{\mathbb{E}\left[\mathbb{1}_{S_i}(\mu) \mid Y = y\right]} + \mathbb{P}(\mu \notin S_i \mid Y = y)f_{\mathrm{loc},\sigma^2}(y)$$

$$= \mathbb{E}[\mathbb{1}_{S_i}(\mu)\mu \mid Y = y] + \mathbb{E}[\mathbb{1}_{S_i^c}(\mu)\mathbb{E}_{\mu' \sim Q_0|_{S_i}}[\mu' \mid Y = y] \mid Y = y]$$

Hence, using the fact that $\mathbb{P}_{(\mu,y)}(\mu \notin S_i) \le \varepsilon$, we have

$$\left\|f_{\sigma^2} - f_{\mathrm{loc},\sigma^2}\right\|_{L^2(P)}^2 = \mathbb{E}_y\left|\mathbb{E}_{\mu \sim Q_0}\left[\mathbb{1}_{S_i^c}(\mu)\left(\mu - \mathbb{E}_{\mu' \sim Q_0|_{S_i}}[\mu' \mid Y = y]\right)\Big|Y = y\right]\right|^2$$

$$\le \mathbb{E}_{(\mu,y)}\left[\mathbb{1}_{S_i^c}(\mu)\left|\mu - \mathbb{E}_{\mu' \sim Q_0|_{S_i}}[\mu' \mid Y = y]\right|^2\right]$$

$$\le 2\mathbb{E}_{(\mu,y)}\left[\mathbb{1}_{S_i^c}(\mu)\left|\mu - \widehat{\mu}_i\right|^2 + \mathbb{1}_{S_i^c}(\mu)\left|\widehat{\mu}_i - \mathbb{E}_{\mu' \sim Q_0|_{S_i}}[\mu' \mid Y = y]\right|^2\right]$$

$$\le 2\mathbb{E}_{(\mu,y)}\left[\mathbb{1}_{S_i^c}(\mu)\left|\mu - \widehat{\mu}_i\right|^2\right] + 2\varepsilon R'^2$$

$$\le 2\int_0^\infty \mathbb{P}_{(\mu,y)}\left(|\mu - \widehat{\mu}_i|^2 \ge 2R'^2 + a\right)da + 6\varepsilon R'^2$$

$$\le 2\int_0^\infty \mathbb{P}_{(\mu,y)}\left(|\mu - \widehat{\mu}_i| \ge R' + \sqrt{\frac{a}{2}}\right)da + 6\varepsilon R'^2. \tag{25}$$

But using Lemma 5.4, we have

$$\mathbb{P}\left(|\mu - \widehat{\mu}_i| \ge 3R + r\right) \le k'e^{-\frac{r^2}{8\sigma^2}},$$

which implies

$$\mathbb{P}_{(\mu,y)}\left(|\mu - \widehat{\mu}_i| \ge R' + \sqrt{\frac{a}{2}}\right) \le k'e^{-\frac{\left(2\sqrt{2}\sigma\sqrt{\ln\left(\frac{k'}{\varepsilon}\right)}+\sqrt{\frac{a}{2}}\right)^2}{8\sigma^2}} \le k'e^{-\frac{8\sigma^2\ln\left(\frac{k'}{\varepsilon}\right)+a/2}{8\sigma^2}} = \varepsilon e^{-\frac{a}{16\sigma^2}}. \tag{26}$$

Plugging Equation (26) into Equation (25),

$$\left\|f_{\sigma^2} - f_{\mathrm{loc},\sigma^2}\right\|_{L^2(P)}^2 \le 2\int_0^\infty \varepsilon e^{-\frac{a}{16\sigma^2}}\,da + 6\varepsilon R'^2 = 32\varepsilon\sigma^2 + 6\varepsilon R'^2.$$

$\square$

## 5.2 Approximation within a Voronoi cell

**Lemma 5.5.** *Let $\mathcal{C} = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_{k'}\}$ and $V_1, \ldots, V_{k'}$ be the corresponding Voronoi partition. Suppose that $\mathcal{C}$ is a complete set of $R$-warm starts. Then*

$$\int_{V_i} \left( \frac{dP_{V_i}}{d\gamma_{\widehat{\mu}_i,1}} \right)^{p+1} d\gamma_{\widehat{\mu}_i,1} \leq \exp \left( \frac{pR^2}{2} \right). \tag{27}$$

*Proof.* We can write

$$P_{V_i} = Q_0(V_i) \frac{P^{V_i}}{Q_0(V_i)} + Q_0(V_i^c) \frac{P^{V_i^c}}{Q_0(V_i^c)};$$

note that $Q_0(V_i)$ and $Q_0(V_i^c)$ are the normalizing constants for the respective probability measures. By convexity, it suffices to show that (27) holds for $P^{V_i}$ and $P^{V_i^c}$ in place of $P_{V_i}$.

For $P^{V_i}$, (27) follows directly from Lemma A.4. For $P^{V_i^c}$, noting that it is a convex mixture of $\mathcal{N}(\mu_j, I_n)$ for $\mu_j \notin V_i$, again by convexity it suffices to show that (27) holds for one such $\mathcal{N}(\mu_j, I_n)$. Suppose that $\mu_j \in V_j$, $j \neq i$. Note that by definition of the Voronoi cell, $|\mu_j - \widehat{\mu}_i| \geq |\mu_j - \widehat{\mu}_j|$ so $\frac{d\gamma_{\widehat{\mu}_j,1}}{d\gamma_{\widehat{\mu}_i,1}} \leq 1$ on $V_i$. Then

$$\int_{V_i} \left( \frac{d\gamma_{\mu_j,1}}{d\gamma_{\widehat{\mu}_i,1}} \right)^{p+1} d\gamma_{\widehat{\mu}_i,1} = \int_{V_i} \left( \frac{d\gamma_{\mu_j,1}}{d\gamma_{\widehat{\mu}_i,1}} \right)^p d\gamma_{\mu_j,1}$$

$$\leq \int_{V_i} \left( \frac{d\gamma_{\mu_j,1}}{d\gamma_{\widehat{\mu}_j,1}} \right)^p d\gamma_{\mu_j,1}$$

$$= \int_{V_i} \left( \frac{d\gamma_{\mu_j,1}}{d\gamma_{\widehat{\mu}_j,1}} \right)^{p+1} d\gamma_{\widehat{\mu}_j,1} \leq e^{\frac{aR^2}{2}} \tag{28}$$

where the last step follows from Lemma A.4, as $\|\mu_j - \widehat{\mu}_j\| \leq R$. $\square$

**Lemma 5.6** (Existence of low-degree polynomial). *Suppose that $f = f_{\sigma^2}$ as in (7), $Q_0$ is supported on $B_D$, $\mathcal{C} = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_{k'}\}$ is a complete set of $R$-warm starts for $Q_0$, and let $V_1, \ldots, V_{k'}$ be the corresponding Voronoi partition. For each $i$, there is a polynomial $g_i$ of degree at most $O\left( \left( \frac{R}{\sigma} + \sqrt{\ln \left( \frac{k'}{\varepsilon} \right)} \right)^6 \ln \left( \frac{1}{\varepsilon} \right)^4 \right)$ such that*

$$\|f - g_i\|_{P_{V_i}}^2 \lesssim \varepsilon^2 \left( R^2 + \sigma^2 \ln \left( \frac{k'}{\varepsilon} \right) \right) \quad \text{and} \quad \|g_i\|_{\gamma_{\widehat{\mu}_i,\sigma^2}} \leq D.$$

*Proof.* Let $R' = 3R + 2\sqrt{2}\sigma \sqrt{\ln \left( \frac{k'}{\varepsilon} \right)}$. First, note that from Lemma 5.3,

$$\sum_{i=1}^{k'} \|f_{\sigma^2} - f_{S_i,\sigma^2}\|_{L^2(P_{V_i})}^2 = \|f_{\sigma^2} - f_{\text{loc},\sigma^2}\|_{L^2(P)}^2 \lesssim \left( R^2 + \sigma^2 \ln \left( \frac{k'}{\varepsilon} \right) \right) \varepsilon^2, \tag{29}$$

where $S_i$ is the expanded neighborhood in Lemma 5.3. By Lemma 4.5 applied to $P_{V_i}$, $f_{\text{loc},\sigma^2}$ is $(O(R'/\sigma), \sigma)$-Gaussian-noise-sensitive. By Lemma 5.5, for each $p > 0$, $\left\| \frac{dP_{V_i}}{d\gamma_{\widehat{\mu}_i,\sigma^2}} \right\|_{L^{1+p}(\gamma_{\widehat{\mu}_i,\sigma^2})} \leq$

$e^{\frac{a(R/\sigma)^2}{2}}$. Hence $P_{V_i}$ satisfies the conditions of Lemma 4.5, and there exists a polynomial $g_i$ of degree at most $O\left(\left(\frac{R'}{\sigma}\right)^4 \ln\left(\frac{1}{\varepsilon}\right)^4 \max\left\{\left(\frac{R'}{\sigma}\right)^2, \ln\left(\frac{1}{\varepsilon}\right)\right\}\right) = O\left(\left(\frac{R'}{\sigma}\right)^6 \ln\left(\frac{1}{\varepsilon}\right)^4\right)$ such that

$$\left\|f_{\mathrm{loc},\sigma^2} - g_i\right\|^2_{P_{V_i}} = \left\|f^{S_i} - g_i\right\|^2_{P_{V_i}} \le \varepsilon^2 R'^2. \tag{30}$$

Combining (29) and (30) gives the result. Finally, from Lemma 4.5,

$$\|g_i\|_{\gamma_{\widehat{\mu}_i,\sigma^2}} \le \left\|f_{\mathrm{loc},\sigma^2}\right\|_{P_{V_i}} \le D.$$

$\square$

## 5.3   Piecewise polynomial regression

For $(x_i, y_i)$ generated according to (5), by (9) we wish to find the least-squares solution to

$$y_i + ts_t(y_i) \approx x_i.$$

where by (7) we parameterize

$$f_{\sigma^2}(y) = y + \sigma^2 s_t(y) = \sum_{|\mathbf{k}| \le d} b_{\mathbf{k}} h_{\mathbf{k}}(y - \widehat{\mu}),$$

where $\widehat{\mu}$ is a warm start. Equivalently, we wish to find the least-squares solution to

$$\sum_{|\mathbf{k}| \le d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i - \widehat{\mu}) \approx y_i + \frac{\sigma^2}{t}(x_i - y_i) = \left(1 - \frac{\sigma^2}{t}\right) y_i + \frac{\sigma^2}{t} x_i.$$

We can restrict to $(b_{\mathbf{k}})_{|\mathbf{k}| \le d} \in B_R$ (in $\mathbb{R}^{n\binom{[n]}{\le d}}$), so we solve

$$\underset{(b_{\mathbf{k}})_{|\mathbf{k}| \le d} \in B_R}{\arg\min} \widehat{L}((b_{\mathbf{k}})_{|\mathbf{k}| \le d}) \quad \text{where} \quad \widehat{L}((b_{\mathbf{k}})_{|\mathbf{k}| \le d}) = \frac{1}{N} \sum_{i=1}^{N} \left| \sum_{|\mathbf{k}| \le d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i - \widehat{\mu}) - \left(\left(1 - \frac{\sigma^2}{t}\right) y_i + \frac{\sigma^2}{t} x_i\right) \right|^2.$$

Let $\overline{x}_i = \mathbb{E}[X|Y = y_i] = y_i + t\nabla V_t(y_i)$ and $x_i = \overline{x}_i + \zeta_i$, so that conditioned on $y_i$, $\zeta_i$ is mean-zero noise. Let $\eta_i = \frac{\sigma^2}{t}\zeta_i$. Let $z_i = \mathbb{E}[\mu|Y = y_i] = \left(1 - \frac{\sigma^2}{t}\right) y_i + \frac{\sigma^2}{t}\overline{x}_i$. Then, for each Voronoi cell $V_j$, we calculate the empirical loss for samples where $y_i$ falls into $V_j$,

$$\widehat{L}^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}| \le d}) = \frac{1}{N} \sum_{i:y_i \in V_j} \left| \sum_{|\mathbf{k}| \le d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i - \widehat{\mu}_j) - (z_i + \eta_i) \right|^2,$$

and the empirical risk minimizer (ERM) for the $j$th Voronoi cell as

$$(\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \le d} := \underset{(b_{\mathbf{k}})_{|\mathbf{k}| \le d} \in B_M}{\arg\min} \widehat{L}^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}| \le d}) \quad . \tag{31}$$

We then combine these ERM solutions for all Voronoi cells and define the following piecewise polynomial function to approximate the score:

$$\widehat{g}(y) = \sum_{j=1}^{k'} \mathbb{1}_{V_j}(y) \sum_{|\mathbf{k}| \leq d} \widehat{b}_{\mathbf{k}}^{(j)} h_{\mathbf{k}}(y - \widehat{\mu}_j), \tag{32}$$

where $k'$ is the total number of Voronoi cells. To show that $\widehat{g}$ is a decent approximation for the score function $f_{\sigma^2}$, we need to bound the generalization error of these ERM solutions (31). To do this, we derive moment bounds for the features and the noise. Let

$$L^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}| \leq d}) = \mathbb{E} \left| \sum_{|\mathbf{k}| \leq d} \mathbb{1}_{V_j}(Y) \left( b_{\mathbf{k}} h_{\mathbf{k}}(Y - \widehat{\mu}_j) - (Z + \eta) \right) \right|^2$$

denote the population risk in $V_j$. We further define the shifted error of the empirical loss with respect to the population loss as

$$\text{err}^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}| \leq d}) := (\widehat{L}^{(j)} - L^{(j)})((b_{\mathbf{k}})_{|\mathbf{k}| \leq d}) + \frac{1}{N} \sum_{i:y_i \in V_j} \left\langle \sum_{|\mathbf{k}| \leq d} z_i, \eta_i \right\rangle - \frac{1}{N} \sum_{i:y_i \in V_j} \left( |\eta_i|^2 - \mathbb{E}|\eta_i|^2 \right)$$

$$= \frac{1}{N} \sum_{i:y_i \in V_j} \left\langle \sum_{|\mathbf{k}| \leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i - \widehat{\mu}_j), \eta_i \right\rangle$$

$$+ \frac{1}{N} \sum_{i:y_i \in V_j} \left( \left| \sum_{|\mathbf{k}| \leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i - \widehat{\mu}_j) - z_i \right|^2 - \mathbb{E} \left| \sum_{|\mathbf{k}| \leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i - \widehat{\mu}_j) - z_i \right|^2 \right). \tag{33}$$

The following lemma is standard and relates the generalization gap of ERM to $\text{err}^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}| \leq d})$. As we see shortly, the terms in the definition (33) that do not depend on $(b_{\mathbf{k}})_{|\mathbf{k}| \leq d}$ cancel out in the calculation of the generalization gap.

**Lemma 5.7** (Uniform convergence $\Rightarrow$ Generalization gap). *The generalization gap for the ERM solution can be bounded as*

$$L^{(j)} \left( (\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \leq d} \right) - \min_{(b_{\mathbf{k}})_{|\mathbf{k}| \leq d} \in B_M} L^{(j)} \left( (b_{\mathbf{k}})_{|\mathbf{k}| \leq d} \right) \leq 2 \max_{(b_{\mathbf{k}})_{|\mathbf{k}| \leq d} \in B_M} \left| \text{err}((b_{\mathbf{k}})_{|\mathbf{k}| \leq d}) \right|.$$

*Proof.* To bound the generalization gap of the ERM solution $(\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \leq d_\ell, 1 \leq j \leq n}$ output by Algorithm 1, we compare its loss with the loss of an arbitrary point $(\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}| \leq d}$ in $B_L(0)$:

$$L^{(j)} \left( (\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \leq d} \right) - L^{(j)} \left( (\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}| \leq d} \right) = L^{(j)} \left( (\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \leq d} \right) - \widehat{L}^{(j)} \left( (\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \leq d} \right)$$

$$+ \widehat{L}^{(j)} \left( (\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}| \leq d} \right) - \widehat{L}^{(j)} \left( (\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}| \leq d} \right)$$

$$+ \widehat{L}^{(j)} \left( (\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}| \leq d} \right) - L^{(j)} \left( (\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}| \leq d} \right). \tag{34}$$

Note that $\widehat{L}^{(j)}\left((\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}|\leq d}\right) - \widehat{L}^{(j)}\left((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}\right) \leq 0$ by the definition of $(\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}|\leq d}$ as a minimizer of $\widehat{L}^{(j)}$. Furthermore, by (33),

$$L^{(j)}\left((\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}|\leq d}\right) - \widehat{L}^{(j)}\left((\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}|\leq d}\right) + \widehat{L}^{(j)}\left((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}\right) - L\left((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}\right)$$

$$= \mathrm{err}^{(j)}((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}) - \mathrm{err}^{(j)}((\widehat{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}|\leq d}) \leq 2 \max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M(0)} \left|\mathrm{err}^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}|\leq d})\right|. \tag{35}$$

$\square$

Based on Lemma 5.7 and Equation (33), to bound the generalization gap it suffices to bound

$$\max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M} \left|\mathrm{err}((b_{\mathbf{k}})_{|\mathbf{k}|\leq d})\right| \leq \max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M} \frac{1}{N} \sum_{i:y_i\in V_j} \left\langle \sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}}h_{\mathbf{k}}(y_i - \widehat{\mu}_j), \eta_i \right\rangle$$

$$+ \max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M} \frac{1}{N} \sum_{i:y_i\in V_j} \left( \left|\sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}}h_{\mathbf{k}}(y_i - \widehat{\mu}_j) - z_i\right|^2 - \mathbb{E}\left|\sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}}h_{\mathbf{k}}(y_i - \widehat{\mu}_j) - z_i\right|^2 \right). \tag{36}$$

Let $h(y) \in \mathbb{R}^{\binom{[n]}{\leq d}}$ be given by $h(y)_{\mathbf{k}} = h_{\mathbf{k}}(y)$. Let $B \in \mathbb{R}^{n\times\binom{[n]}{\leq d}}$ with columns $b_{\mathbf{k}}$. Then $\left|(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\right| = \|B\|_F$ so

$$\max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M} \sum_{i:y_i\in V_j} \sum_{|\mathbf{k}|\leq d} \langle b_{\mathbf{k}}h_{\mathbf{k}}(y_i - \widehat{\mu}_j), \eta_i \rangle = \max_{\|B\|_F\leq M} \mathrm{Tr}\left( B \sum_{i:y_i\in V_j} h(y_i - \widehat{\mu}_j)\eta_i^\top \right)$$

$$= M\left\| \sum_{i:y_i\in V_j} h(y_i - \widehat{\mu}_j)\eta_i^\top \right\|_F. \tag{37}$$

We carefully bound the moments of each of the terms in Equation (36) in Lemmas 5.9 and 5.11, leading to high-probability bounds given in Lemmas 5.10 and 5.12, respectively. From this, we obtain the following high probability bound on the generalization gap.

**Lemma 5.8** (Generalization gap). *Suppose $Q_0$ is supported on $B_D(0)$ and $\mathcal{C}$ is a complete set of $R$-warm starts for $Q_0$ with $|\mathcal{C}'| = k'$ and $d \geq (R/\sigma)^2$. Given $N$ pairs of samples $(\mu_i, y_i)$ distributed as $y_i = \mu_i + \xi_i$, for $\mu_i \sim Q_0$ and $\xi_i \sim \mathcal{N}(0, \sigma^2 I_n)$ with*

$$N = \Omega\left( \frac{n\left(k'(1+M+D)^2\right)^2 \left(\ln(k/\delta)\sqrt{n/d+1}\right)^{4(d+1)}}{\varepsilon^4} \right),$$

*we have with probability at least $1 - \delta$ that*

$$\max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M(0)} \left|\mathrm{err}^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}|\leq d})\right| \leq \frac{\varepsilon^2}{k'}.$$

*Furthermore, the piecewise polynomial approximation of the score defined in Equation (32) satisfies*

$$\|\widehat{g} - f\|_{L^2(P)}^2 - \|\widetilde{g} - f\|_{L^2(P)}^2 \leq \varepsilon^2.$$

*for any other piece-wise polynomial $\widetilde{g}$ on the Voronoi cells whose coefficients satisfy $(\widetilde{b}_{\mathbf{k}}^{(j)})_{|\mathbf{k}|\leq d} \in B_M$ for each cell $j$.*

31

We defer the proof to Section 5.3.2.

### 5.3.1 Moment and high-probability bounds

For the following lemmas, we assume the following: $Q_0$ is supported on $B_D(0)$, $\mathcal{C}$ is a complete set of $R$-warm starts for $Q_0$, and $V_1, \ldots, V_{k'}$ is the corresponding Voronoi partition. We are given iid samples $(\mu_i, y_i), i = 1, \ldots, N$, where $y_i = \mu_i + \xi_i$, $\mu_i \sim Q_0$, and $\xi_i \sim \mathcal{N}(0, \sigma^2 I_n)$. Recall that $h_{\mathbf{k}}$ are the $\sigma$-rescaled Hermite polynomials and $h(y) = (h_{\mathbf{k}}(y))_{|\mathbf{k}| \leq d}$. We moreover fix the Voronoi cell $V_j$ corresponding to $\widehat{\mu}_j$, and (to simplify notation) suppose that $\widehat{\mu}_j = 0$.

**Lemma 5.9** (Moment bound for first term). *We have the following uniform convergence bound for $p$ even:*

$$\left\| \left\| \sum_{i=1}^{N} \mathbb{1}_{V_j}(y_i) h(y_i) \eta_i^\top \right\|_F \right\|_{2p} = O\left( Dn^{1/2} \left( \frac{n}{d} + 1 \right)^{d/2} \sqrt{p} \left( e(p-1) \right)^d \sqrt{N} \right).$$

We will use properties of sub-exponential random variables; see Appendix B for background.

*Proof.* In the following, the expectations without index refer to the mixture measure. Let $\mu_i^{(2)}$ be independent and identically distributed to $\mu_i | Y = y_i$. Without loss of generality suppose the samples in which $y_i$ fall into $V_j$ are exactly $y_1, \ldots, y_{N_j}$. Using the symmetrization technique with Jensen's inequality:

$$\mathbb{E} \left\| \sum_{i=1}^{N_j} h(y_i) \eta_i^\top \right\|_F^{2p} = \mathbb{E} \left\| \sum_{i=1}^{N_j} h(y_i)(\mu_i - \mathbb{E}[\mu_i | Y = y_i])^\top \right\|_F^{2p} \tag{38}$$

$$= \mathbb{E} \left\| \sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mathbb{E}[\mu_i | Y = y_i])^\top \right\|_F^{2p} \tag{39}$$

$$= \mathbb{E} \left\| \sum_{i=1}^{N_j} h(y_i)(\mathbb{E}_{\mu_i}[(\mu_i^{(2)} - \mu_i)|Y = y_i])^\top \right\|_F^{2p} \tag{40}$$

$$\leq \mathbb{E} \left\| \sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mu_i)^\top \right\|_F^{2p}. \tag{41}$$

Note that conditioned on any fixed value of $y_i$, the variable $\mu_i^{(2)} - \mu_i$ is the difference between two iid random variables, each of whose norm is bounded by $D$. Therefore, given a fixed $\{y_i\}_{i=1}^{N_j}$, the variable $\left\| \sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mu_i)^\top \right\|_F^2$ is a quadratic form of a sub-Gaussian vector. From the properties of sub-Gaussian variables, for every $1 \leq k \leq n$ and $|\mathbf{k}| \leq d$,

$$\left( \sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)(\mu_i^{(2)} - \mu_i)_k \right)^2$$

is sub-exponential with parameter $\left(O\left(\left(\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)^2 D^2\right)^2\right), O\left(\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)^2 D^2\right)\right)$. Therefore, noting we do not have independence when we enumerate the variables over $\mathbf{k}$, by Proposition B.4, $\sum_{\mathbf{k}}(\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)(\mu_i^{(2)} - \mu_i)_k)^2$ is sub-exponential with parameters

$$\left(\binom{n+d}{d} O\left(\sum_{|\mathbf{k}|\leq d}\left(\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)^2 D^2\right)^2\right), \binom{n+d}{d} O\left(\max_{|\mathbf{k}|\leq d}\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)^2 D^2\right)\right).$$

Finally summing over $1 \leq k \leq n$, from the independence of the coordinates of $(\mu^{(2)} - \mu)$, by Proposition B.3, the variable

$$X := \left\|\sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mu_i)^\top\right\|_F^2 = \sum_{k=1}^n \sum_{|\mathbf{k}|\leq d}\left(\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)((\mu_i^{(2)} - \mu_i)^\top)_k\right)^2$$

is sub-exponential with parameters

$$(v^2, \alpha) := \left(n\binom{n+d}{d} O\left(\sum_{|\mathbf{k}|\leq d}\left(\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)^2 D^2\right)^2\right), n\binom{n+d}{d} O\left(\max_{|\mathbf{k}|\leq d}\sum_{i=1}^{N_j} h_{\mathbf{k}}(y_i)^2 D^2\right)\right)$$

and has expectation

$$\mathbb{E}\left\|\sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mu_i)^\top\right\|_F^2 = \sum_{i=1}^{N_j} |h(y_i)|^2 \, \mathbb{E}\left|\mu_i - \mu_i^{(2)}\right|^2 = 4D^2 \sum_{i=1}^{N_j} |h(y_i)|^2.$$

But this implies

$$\nVdash P\left(X - \mathbb{E}[X] \geq t\right) \leq e^{-t/\kappa},$$

for parameter

$$\kappa = O\left(\max(v, \alpha)\right) = O\left(D^2 n\binom{n+d}{d}\sum_{i=1}^{N_j} |h(y_i)|^2\right).$$

Thus, using the second property for sub-exponential variables in Proposition B.2, for some constant $c_1$ we have

$$\mathbb{E}\left[\left\|\sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mu_i)^\top\right\|_F^{2p}\left|\{y_i\}_{i=1}^{N_j}\right.\right] \leq O(p\kappa + \mathbb{E}X)^p = \left(c_1 p D^2 n\binom{n+d}{d}\sum_{i=1}^{N_j} |h(y_i)|^2\right)^p.$$

$$(42)$$

Then by Lemma A.3,

$$\mathbb{E}\left\|\sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mu_i)^\top\right\|_F^{2p} = \mathbb{E}\left[\mathbb{E}\left[\left\|\sum_{i=1}^{N_j} h(y_i)(\mu_i^{(2)} - \mu_i)^\top\right\|_F^{2p}\left|\{y_i\}_{i=1}^{N_j}\right.\right]\right]$$

$$\leq \mathbb{E}\left(c_1 p D^2 n\binom{n+d}{d}\sum_{i=1}^{N_j} |h(y_i)|^2\right)^p$$

Therefore, given $\gamma = \gamma_{\widehat{\mu}_j, \sigma^2}$ is the Gaussian around the warm start point $\widehat{\mu}_j$, and Lemma A.3 applied to $|h|^p$ (which is a polynomial as $p$ is even) and Lemma A.2,

$$\left\| \left\| \left\| \sum_{i=1}^{N} \mathbb{1}_{V_j}(y_i) h(y_i) \eta_i^\top \right\| \right\|_F \right\|_{2p} \leq \sqrt{c_1} \sqrt{p} D \left( n \binom{n+d}{d} \right)^{1/2} \sqrt{\left\| \left\| \sum_{i=1}^{N_j} |h(y_i)|^2 \right\| \right\|_p}$$

$$\leq \sqrt{c_1} \sqrt{p} D \left( n \binom{n+d}{d} \right)^{1/2} \sqrt{\sum_{i=1}^{N_j} \left\| |h(y_i)|^2 \right\|_p}$$

$$\leq \sqrt{c_1} \sqrt{p} D \left( n \binom{n+d}{d} \right)^{1/2} \sqrt{\sum_{i=1}^{N_j} e^{2\sqrt{d/p}(R/\sigma)} \left\| |h(y_i)|^2 \right\|_{L^p(\gamma)}}$$

$$\leq \sqrt{c_1} \sqrt{p} D \left( n \binom{n+d}{d} \right)^{1/2} \sqrt{\sum_{i=1}^{N_j} e^{2\sqrt{d/p}(R/\sigma)} (p-1)^{2d} \left\| |h(y_i)|^2 \right\|_{L^2(\gamma)}}$$

$$\leq c_2 \sqrt{p} D n^{1/2} \left( \frac{n}{d} + 1 \right)^{d/2} (e(p-1))^d \sqrt{N},$$

for some universal constant $c_2$, and we used the upper bound $\binom{n}{k} \leq \frac{1}{e} \left( \frac{en}{k} \right)^k$. $\qquad \square$

**Lemma 5.10** (High probability bound for first term). *We have*

$$\mathbb{P}\left( \left\| \sum_{i=1}^{N} \mathbb{1}_{V_j}(y_i) h(y_i) \eta_i^\top \right\|_F \geq t \right) \leq \exp\left( -\Omega\left( \left( \frac{n}{d} + 1 \right)^{-1/2} \left( \frac{t}{D\sqrt{N} n^{1/2}} \right)^{1/(d+1)} \right) \right).$$

*Proof.* By Lemma 5.9 and Markov's inequality,

$$\mathbb{P}\left( \left\| \sum_{i=1}^{N} \mathbb{1}_{V_j}(y_i) h(y_i) \eta_i^\top \right\|_F \geq t \right) = \mathbb{P}\left( \left\| \sum_{i=1}^{N} \mathbb{1}_{V_j}(y_i) h(y_i) \eta_i^\top \right\|_F^{2p} \geq t^{2p} \right)$$

$$\leq \frac{\mathbb{E} \left\| \sum_{i=1}^{N} \mathbb{1}_{V_j}(y_i) h(y_i) \eta_i^\top \right\|_F^{2p}}{t^{2p}}$$

$$\leq \left( \frac{D n^{1/2} \left( \frac{n}{d} + 1 \right)^{d/2} (e(p-1))^{d+1} \sqrt{N}}{t} \right)^{2p}.$$

Picking

$$p - 1 = \Theta\left( \frac{\left( t/(D\sqrt{N} n^{1/2}) \right)^{1/(d+1)}}{\left( \frac{n}{d} + 1 \right)^{1/2}} \right)$$

with appropriate constant, we have

$$\mathbb{P}\left( \left\| \sum_{i=1}^{N} \mathbb{1}_{V_j}(y_i) h(y_i) \eta_i^\top \right\|_F \geq t \right) \leq \exp\left( -\Omega\left( \left( \frac{n}{d} + 1 \right)^{-1/2} \left( \frac{t}{D\sqrt{N} n^{1/2}} \right)^{1/(d+1)} \right) \right).$$

$\qquad \square$

**Lemma 5.11** (Moment bound for second term)**.** *Assume $d \geq (R/\sigma)^2$. We have the following uniform convergence bound:*

$$\left\| \max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M} \sum_{i=1}^N \mathbb{1}_{V_j}(y_i)\left( \left| \sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i) - z_i \right|^2 - \mathbb{E}\left| \sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i) - z_i \right|^2 \right) \right\|_{2p}$$

$$= O\left( M(M+D)n^{1/2}e^d \left( \frac{n}{d}+1 \right)^d p\,(4(4p-1))^{2d}\sqrt{N} + \sqrt{pN}D^2 \right).$$

*Proof.* Without loss of generality, suppose the first $N_j$ samples $y_1, \ldots, y_{N_j}$ are in $V_j$. Again using symmetrization for the $y_i$'s,

$$\left\| \max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M} \sum_{i=1}^{N_j}\left( \left| \sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i) - z_i \right|^2 - \mathbb{E}\left| \sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i) - z_i \right|^2 \right) \right\|_{2p}.$$

$$\leq \left\| \max_{(b_{\mathbf{k}})_{|\mathbf{k}|\leq d}\in B_M} \sum_{i=1}^{N_j}\left( \left| \sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i^{(2)}) - z_i^{(2)} \right|^2 - \mathbb{E}\left| \sum_{|\mathbf{k}|\leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i) - z_i \right|^2 \right) \right\|_{2p}$$

$$\leq \left\| \max_{\|B\|_F\leq M} \mathrm{tr}\left( B \left( \sum_{i=1}^{N_j} h(y_i)h(y_i)^\top - \sum_{i=1}^{N_j} h(y_i^{(2)})h(y_i^{(2)})^\top \right) B^\top \right) \right\|_{2p}$$

$$+ \left\| \max_{\|B\|_F\leq M} \left\{ \sum_{i=1}^{N_j} z_i^\top B h(y_i) - \sum_{i=1}^{N_j} z_i^{(2)\top} B h(y_i^{(2)}) \right\} \right\|_{2p} + \left\| \sum_{i=1}^{N_j}\left( \left| z_i^{(2)} \right|^2 - |z_i|^2 \right) \right\|_{2p}$$

$$\leq \left\| M^2 \left\| \sum_{i=1}^{N_j} h(y_i)h(y_i)^\top - \sum_{i=1}^{N_j} h(y_i^{(2)})h(y_i^{(2)})^\top \right\|_{\mathrm{op}} \right\|_{2p}$$

$$+ \left\| \max_{\|B\|_F\leq M} \left\langle B, \sum_{i=1}^{N_j} z_i h(y_i)^\top - \sum_{i=1}^{N_j} z_i^{(2)} h(y_i^{(2)})^\top \right\rangle \right\|_{2p} + \left\| \sum_{i=1}^{N_j}\left( \left| z_i^{(2)} \right|^2 - |z_i|^2 \right) \right\|_{2p}$$

$$\leq M^2 \left\| \left\| \sum_{i=1}^{N_j} h(y_i)h(y_i)^\top - \sum_{i=1}^{N_j} h(y_i^{(2)})h(y_i^{(2)})^\top \right\|_F \right\|_{2p}$$

$$+ M \left\| \left\| \sum_{i=1}^{N_j} z_i h(y_i)^\top - \sum_{i=1}^{N_j} z_i^{(2)} h(y_i^{(2)})^\top \right\|_F \right\|_{2p} + \left\| \sum_{i=1}^{N_j}\left( \left| z_i^{(2)} \right|^2 - |z_i|^2 \right) \right\|_{2p}.$$

For the third term, note that for each $1 \leq i \leq N$:

$$\left| \left| z_i^{(2)} \right|^2 - |z_i|^2 \right| \leq 2D^2,$$

and this has mean zero. Hence the third term is sub-Gaussian with parameter $O(D^2\sqrt{N})$. Therefore

$$\left\| \sum_{i=1}^{N_j}\left( \left| z_i^{(2)} \right|^2 - |z_i|^2 \right) \right\|_{2p} \leq O\left( \sqrt{pN}D^2 \right). \tag{43}$$

35

For the second term, note that for all $1 \leq i \leq N_j$, $1 \leq k \leq n$, and $|\mathbf{k}| \leq \binom{n}{d}$, given $\gamma = \gamma_{\widehat{\mu}_j, \sigma^2}$ is the Gaussian around the warm start point $\widehat{\mu}_j$: by Cauchy-Schwarz, Lemma A.3 applied to $h_{\mathbf{k}}^p$, Lemma A.2, and the assumption that $d \geq (R/\sigma)^2$,

$$
\begin{aligned}
\left\| z_{ik} h_{\mathbf{k}}(y_i) - z_{i_k}^{(2)} h_{\mathbf{k}}(y_i^{(2)}) \right\|_p &\leq 2 \left\| z_{ik} h_{\mathbf{k}}(y_i) \right\|_p \\
&\leq 2 \left\| z_{ik} \right\|_{2p} \left\| h_{\mathbf{k}}(y_i) \right\|_{2p} \\
&\leq 2D \cdot 2^{\sqrt{d/(2p)}(R/\sigma)} \left\| h_{\mathbf{k}}(y_i) \right\|_{L^{2p}(\gamma)} \\
&\leq 2D \cdot 2^{\sqrt{d/(2p)}(R/\sigma)} (2p-1)^d \left\| h_{\mathbf{k}}(y_i) \right\|_{L^2(\gamma)} \\
&\leq 2D \cdot (2(2p-1))^d .
\end{aligned}
$$

Therefore, from the Rosenthal inequality [IS01] with constant given by Pinelis [NP78], there is a universal constant $c$ such that

$$
\mathbb{E} \left| \sum_{i=1}^{N_j} z_{ik} h_{\mathbf{k}}(y_i) - z_{i_k}^{(2)} h_{\mathbf{k}}(y_i^{(2)}) \right|^p
$$

$$
\leq (cp)^p \max \left\{ \sum_{i=1}^{N_j} \mathbb{E} \left| z_{ik} h_{\mathbf{k}}(y_i) - z_{i_k}^{(2)} h_{\mathbf{k}}(y_i^{(2)}) \right|^p , \left( \sum_{i=1}^{N_j} \mathbb{E} \left| z_{ik} h_{\mathbf{k}}(y_i) - z_{i_k}^{(2)} h_{\mathbf{k}}(y_i^{(2)}) \right|^2 \right)^{p/2} \right\}
$$

$$
\leq (cp)^p \max \left\{ N_j 2^p D^p (2(2p-1))^{pd}, \left( 4 N_j D^2 6^{2d} \right)^{p/2} \right\}
$$

$$
\leq \left( 2cp \, (6(2p-1))^d \sqrt{N} D \right)^p ,
$$

which according to the norm property of $p$th norm $\left\| \cdot \right\|_p$ implies

$$
\left\| \sum_{|\mathbf{k}| \leq d} \sum_{k=1}^n \left( \sum_{i=1}^{N_j} z_{ik} h_{\mathbf{k}}(y_i) - z_{i_k}^{(2)} h_{\mathbf{k}}^{(2)} \right)^2 \right\|_p \leq \sum_{|\mathbf{k}| \leq d} \sum_{k=1}^n \left\| \left( \sum_{i=1}^{N_j} z_{ik} h_{\mathbf{k}}(y_i) - z_{i_k}^{(2)} h_{\mathbf{k}}^{(2)} \right)^2 \right\|_p
$$

$$
\leq n \binom{n+d}{d} \left( 4cp \, (6(4p-1))^d \sqrt{N} D \right)^2 .
$$

Therefore, using $\binom{n}{k} \leq \frac{1}{e} \left( \frac{en}{k} \right)^k$,

$$
\left\| \left\| \sum_{i=1}^{N_j} z_i h(y_i)^\top - \sum_{i=1}^{N_j} z_i^{(2)} h(y_i^{(2)})^\top \right\|_F \right\|_{2p} \leq n^{1/2} \left( e \left( \frac{n}{d} + 1 \right) \right)^{d/2} 4cp \, (6(4p-1))^d \sqrt{N} D. \qquad (44)
$$

For the first term, for all $|\mathbf{k}_1|, |\mathbf{k}_2| \leq d$, again by Cauchy-Schwarz, Lemma A.3 applied to $h_{\mathbf{k}}^p$,

Lemma A.2, and the assumption that $d \geq (R/\sigma)^2$,

$$
\left\| h_{\mathbf{k}_1}(y_i)h_{\mathbf{k}_2}(y_i) - h_{\mathbf{k}_1}(y_i^{(2)})h_{\mathbf{k}_2}(y_i^{(2)}) \right\|_p \leq 2 \left\| h_{\mathbf{k}_1}(y_i)h_{\mathbf{k}_2}(y_i) \right\|_p
$$
$$
\leq 2 \left\| h_{\mathbf{k}_1}(y_i) \right\|_{2p} \left\| h_{\mathbf{k}_2}(y_i) \right\|_{2p}
$$
$$
\leq 2^{2\sqrt{d/(2p)}(R/\sigma)+1} \left\| h_{\mathbf{k}_1}(y_i) \right\|_{L^{2p}(\gamma)} \left\| h_{\mathbf{k}_2}(y_i) \right\|_{L^{2p}(\gamma)}
$$
$$
\leq 2^{2\sqrt{d/(2p)}(R/\sigma)+1} (2p-1)^{2d} \left\| h_{\mathbf{k}_1}(y_i) \right\|_{L^2(\gamma)} \left\| h_{\mathbf{k}_2}(y_i) \right\|_{L^2(\gamma)}
$$
$$
\leq (4(2p-1))^{2d}.
$$

Therefore, again using Rosenthal's inequality,

$$
\mathbb{E} \left| \sum_{i=1}^{N_j} \left( h_{\mathbf{k}_1}(y_i)h_{\mathbf{k}_2}(y_i) - h_{\mathbf{k}_1}(y_i^{(2)})h_{\mathbf{k}_2}(y_i^{(2)}) \right) \right|^p
$$
$$
\leq (cp)^p \max \left\{ \sum_{i=1}^{N_j} \mathbb{E} \left| h_{\mathbf{k}_1}(y_i)h_{\mathbf{k}_2}(y_i) - h_{\mathbf{k}_1}(y_i^{(2)})h_{\mathbf{k}_2}(y_i^{(2)}) \right|^p, \left( \sum_{i=1}^{N_j} \mathbb{E} \left| h_{\mathbf{k}_1}(y_i)h_{\mathbf{k}_2}(y_i) - h_{\mathbf{k}_1}(y_i^{(2)})h_{\mathbf{k}_2}(y_i^{(2)}) \right|^2 \right)^{p/2} \right\}
$$
$$
\leq \left( cp\, (4(2p-1))^{2d}\, \sqrt{N} \right)^p,
$$

which implies

$$
\left\| \sum_{|\mathbf{k}_1|,|\mathbf{k}_2| \leq d} \left( \sum_{i=1}^{N_j} \left( h_{\mathbf{k}_1}(y_i)h_{\mathbf{k}_2}(y_i) - h_{\mathbf{k}_1}(y_i^{(2)})h_{\mathbf{k}_2}(y_i^{(2)}) \right) \right)^2 \right\|_p
$$
$$
\leq \sum_{|\mathbf{k}_1|,|\mathbf{k}_2| \leq d} \left\| \left( \sum_{i=1}^{N_j} \left( h_{\mathbf{k}_1}(y_i)h_{\mathbf{k}_2}(y_i) - h_{\mathbf{k}_1}(y_i^{(2)})h_{\mathbf{k}_2}(y_i^{(2)}) \right) \right)^2 \right\|_p
$$
$$
\leq \binom{n+d}{d}^2 \left( 2cp\, (4(4p-1))^{2d}\, \sqrt{N} \right)^2.
$$

Therefore, using $\binom{n}{k} \leq \frac{1}{e} \left( \frac{en}{k} \right)^k$,

$$
\left\| \left\| \sum_{i=1}^{N_j} \left( h(y_i)h(y_i)^\top - h(y_i^{(2)})h(y_i^{(2)})^\top \right) \right\|_F \right\|_{2p} \leq \left( e \left( \frac{n}{d}+1 \right) \right)^d 2cp\, (4(4p-1))^{2d}\, \sqrt{N}. \qquad (45)
$$

Combining Equations (45), (44), and (43) completes the proof. $\qquad \square$

**Lemma 5.12** (High-probability bound for second term). *Assume $d \geq (R/\sigma)^2$. We have*

$$
\mathbb{P} \left( \max_{(b_{\mathbf{k}})_{|\mathbf{k}| \leq d} \in B_M(0)} \sum_{i=1}^{N} \left( \left| \sum_{|\mathbf{k}| \leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i) - z_i \right|^2 - \mathbb{E} \left| \sum_{|\mathbf{k}| \leq d} b_{\mathbf{k}} h_{\mathbf{k}}(y_i) - z_i \right|^2 \right) \geq t \right)
$$
$$
\leq \exp \left( -\Omega \left( \left( \frac{n}{d}+1 \right)^{-1/2} \left( \frac{t}{n^{1/2} M(M+D)\sqrt{N}} \right)^{1/(2(d+1))} \wedge \left( \frac{t}{D^2 \sqrt{N}} \right)^2 \right) \right).
$$

*Proof.* We have

$$\mathbb{P}\left(\max_{(b_{\mathbf{k}})_{|\mathbf{k}|\le d}\in B_M(0)}\sum_{i=1}^{N}\left(\left|\sum_{|\mathbf{k}|\le d}b_{\mathbf{k}}h_{\mathbf{k}}(y_i)-z_i\right|^2-\mathbb{E}\left|\sum_{|\mathbf{k}|\le d}b_{\mathbf{k}}h_{\mathbf{k}}(y_i)-z_i\right|^2\right)\ge t\right)$$

$$=\frac{\mathbb{E}\left|\max_{(b_{\mathbf{k}})_{|\mathbf{k}|\le d}\in B_M(0)}\sum_{i=1}^{N}\left(\left|\sum_{|\mathbf{k}|\le d}b_{\mathbf{k}}h_{\mathbf{k}}(y_i)-z_i\right|^2-\mathbb{E}\left|\sum_{|\mathbf{k}|\le d}b_{\mathbf{k}}h_{\mathbf{k}}(y_i)-z_i\right|^2\right)\right|^{2p}}{t^{2p}}$$

$$\le\left(\frac{n^{1/2}M(M+D)e^d\left(\frac{n}{d}+1\right)^d p\,(4(4p-1))^{2d}\sqrt{N}+O(\sqrt{pN}D^2)}{t}\right)^{2p}.$$

Picking

$$p=O\left(\frac{\left(c't/(n^{1/2}M(M+D)\sqrt{N})\right)^{1/(2(d+1))}}{\left(\frac{n}{d}+1\right)^{1/2}}\wedge\left(\frac{c't}{D^2\sqrt{N}}\right)^2\right),$$

for universal constant $c'$ small enough, the proof is complete. $\qquad\square$

### 5.3.2 Generalization gap

*Proof of Lemma 5.8.* Note that from Equation (33), for every fixed Voronoi cell $V_j$:

$$\max_{(b_{\mathbf{k}})_{|\mathbf{k}|\le d}\in B_M}\left|\mathrm{err}^{(j)}((b_{\mathbf{k}})_{|\mathbf{k}|\le d})\right|\le\frac{1}{N}\left\|\sum_{i:y_i\in V_j}h(y_i-\widehat{\mu}_j)\eta_i^{\top}\right\|_F$$

$$+\max_{(b_{\mathbf{k}})_{|\mathbf{k}|\le d}\in B_M}\frac{1}{N}\sum_{i:y_i\in V_j}\left(\left|\sum_{|\mathbf{k}|\le d}b_{\mathbf{k}}h_{\mathbf{k}}(y_i-\widehat{\mu}_j)-z_i\right|^2-\mathbb{E}\left|\sum_{|\mathbf{k}|\le d}b_{\mathbf{k}}h_{\mathbf{k}}(y_i-\widehat{\mu}_j)-z_i\right|^2\right).\quad(46)$$

But from Lemma 5.10, given

$$N=\Omega\left(\frac{k'^2\ln^{2(d+1)}(k'/\delta)D^2n\,(n/d+1)^{(d+1)}}{\varepsilon^4}\right),$$

samples, with probability at least $1-\delta/(2k')$ the absolute value of the first term in Equation (46) is at most $\varepsilon^2/(2k')$. Furthermore, based on Lemma 5.12, given

$$N=\Omega\left(\frac{k'^2\ln^{4(d+1)}(k'/\delta)n(M+D)^4(n/d+1)^{2(d+1)}}{\varepsilon^4}\right)$$

with probability at least $1-\delta/(2k')$ the second term is bounded by $\varepsilon^2/(2k')$. Applying a union bound, the sum of the first and second terms is bounded by $\varepsilon^2$ with probability at least $1-\delta/k'$. This shows the first claim. To show the second claim, note that

$$L^{(j)}((\widehat{b}_{\mathbf{k}})_{|\mathbf{k}|\le d})=\mathbb{E}_{(\mu,Y)}\left|\mathbb{1}_{V_j}(Y)\left(\widehat{g}(Y)-(Z+\eta)\right)\right|^2$$

$$=\mathbb{E}_{(\mu,Y)}\left|\mathbb{1}_{V_j}(Y)\left(\widehat{g}(Y)-Z\right)\right|^2+\mathbb{E}[\mathbb{1}_{V_j}(Y)\,|\eta|^2],$$

where $\widehat{g}$ is defined as in (32) and similarly

$$L^{(j)}((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}) = \mathbb{E}_{(\mu,Y)}\left|\mathbb{1}_{V_j}(Y)\left(\widetilde{g}(Y)-Z\right)\right|^2 + \mathbb{E}[\mathbb{1}_{V_j}(Y)\,|\eta|^2],$$

where $\widetilde{g}$ is defined analogously. Therefore

$$L^{(j)}((\widehat{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}) - L^{(j)}((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}) = \mathbb{E}_{(\mu,Y)}\left|\mathbb{1}_{V_j}(Y)\left(\widehat{g}(Y)-Z\right)\right|^2 - \mathbb{E}_{(\mu,Y)}\left|\mathbb{1}_{V_j}(Y)\left(\widetilde{g}(Y)-Z\right)\right|^2. \tag{47}$$

Summing Equation (47) for $1 \leq j \leq k'$ implies

$$\sum_{j=1}^{k'} L^{(j)}((\widehat{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}) - L^{(j)}((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}) = \|\widehat{g}-f\|^2_{L^2(P)} - \|\widetilde{g}-f\|^2_{L^2(P)}.$$

But from the previous claim and union bound, we know with probability at least $1 - \delta$ each $L^{(j)}((\widehat{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d}) - L^{(j)}((\widetilde{b}_{\mathbf{k}})_{|\mathbf{k}|\leq d})$ is bounded by $\varepsilon^2/k'$. This completes the proof. $\qquad\square$

## 5.4 Maintaining warm starts

We would like to maintain warm starts to the centers of all Gaussians (of non-negligible mass) as we decrease the noise level. By choosing the highest noise level large enough, $\mathbf{0}$ will be a warm start. The key observation is that with high probability, the score function points in a direction close to a mean; this remains true with the estimated score when the error is small.

---

**Algorithm 2** Picking set of warm starts

---

1: **Input:** Score estimate $s$, noise level $\sigma$, data points $y_1, \ldots, y_N \sim Q_{\sigma^2}$ $(N = \Omega\left(\frac{\ln(1/\delta)k}{\alpha_{\min}}\right))$, min weight $\alpha_{\min}$, failure probability $\delta$.
2: For each $i$, let $\widehat{\mu}_i = y_i + \sigma^2 s(y_i)$.
3: Let $U = \{\widehat{\mu}_1, \ldots, \widehat{\mu}_N\}$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Uncovered means
4: Let $\widetilde{R} = C\left(R_0 + \sigma\sqrt{\ln\left(\frac{1}{\alpha_{\min}}\right)}\right)$ for an appropriate constant $C$.
5: **for** $t = 1$ to $C'k\ln\left(\frac{1}{\alpha_{\min}}\right)$ **do** $\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Greedy set cover
6: $\qquad$ Let $\widehat{\mu} = \arg\max_{\mu\in U}|B_\mu(\widetilde{R})\cap U|$.
7: $\qquad$ Let $D = B_{\widehat{\mu}}(\widetilde{R})\cap U$.
8: $\qquad$ $\widehat{\mathcal{C}} \leftarrow \widehat{\mathcal{C}} \cup \{\widehat{\mu}\}$, $U \leftarrow U\backslash D$.
9: **end for**
10: **Output:** Set of warm starts $\widehat{\mathcal{C}}$

---

We first show that we do not lose too much in the score estimate if we only consider the means that are close to the mean that a data point came from.

**Lemma 5.13** (Good score estimation $\implies$ Warm starts). *Consider a mixture of Gaussians satisfying Assumption 1.1 with $\sigma_0^2 = 1$. There is a universal constant such that the following holds. Let $f_{\sigma^2}(y) = y + \sigma^2\nabla\ln q_{\sigma^2}(y)$. Suppose we are given a function $g$ satisfying*

$$\|f - g\|^2_{L^2(Q_{\sigma^2})} \leq (R_0 + \sigma)^2\alpha_{\min}. \tag{48}$$

*Suppose we have $N = \Omega\left(\frac{\ln(1/\delta)k}{\alpha_{\min}}\right)$ samples sampled from $q_{\sigma^2}$. Let $\mathcal{C}$ be the output of Algorithm 2. Then with probability $\geq 1 - \delta$, for radius $\widetilde{R} = C\left(R_0 + \sigma\sqrt{\ln\left(\frac{1}{\alpha_{\min}}\right)}\right)$ with universal constant $C$, the support of $Q_0$ is contained in $\bigcup_{\widehat{\mu}\in\mathcal{C}} B_{\widetilde{R}}(\widehat{\mu})$.*

*Proof.* By (48) and Chebyshev's inequality,

$$Q_{\sigma^2}(|f - g| \geq 4(R_0 + \sigma)) = Q_{\sigma^2}(|f - g|^2 \geq 16(R_0 + \sigma)^2) \leq \frac{\|f - g\|^2_{L^2(q_\sigma^2)}}{16(R_0 + \sigma)^2} \leq \frac{\alpha_{\min}}{16}.$$

Consider drawing $\mu \sim Q_0$, $\xi \sim \mathcal{N}(0, I_n)$, and $Y = \mu + \sigma\xi$. Let $\mathcal{C} = \{\overline{\mu}_1, \ldots, \overline{\mu}_k\}$ be as in Assumption 1.1. Let $f_{\text{loc},\sigma^2} = f_{\text{loc},\sigma^2}^{\mathcal{C},R'}$ be as in Lemma 5.3, where $R' = 3R_0 + 2\sqrt{2}\sigma\sqrt{\ln\left(\frac{k}{\varepsilon'}\right)}$. Then

$$\mathbb{E}\left|f_{\sigma^2}(y) - f_{\text{loc},\sigma^2}(y)\right|^2 \lesssim \left(R_0^2 + \sigma^2\ln\left(\frac{k}{\varepsilon'}\right)\right)\varepsilon'.$$

By choosing $\varepsilon' = \frac{c\alpha_{\min}}{\ln\left(\frac{1}{\alpha_{\min}}\right)}$ for a small enough constant $c$ and noting $\alpha_{\min} \leq \frac{1}{k}$, we obtain that $\mathbb{E}\left|f_{\sigma^2}(y) - f_{\text{loc},\sigma^2}(y)\right|^2 \leq (R_0 + \sigma)^2\alpha_{\min}$. Again by Chebyshev's inequality,

$$\mathbb{P}\left(\left|f(y) - f_{\text{loc},\sigma^2}(y)\right| \geq 4(R_0 + \sigma)\right) \leq \frac{\alpha_{\min}}{16}.$$

Hence

$$\mathbb{P}\left(\left|g(y) - f_{\text{loc},\sigma^2}(y)\right| \geq 8(R_0 + \sigma)\right) \leq \frac{\alpha_{\min}}{8}. \tag{49}$$

Let $V_1, \ldots, V_k$ be the Voronoi partition corresponding to $\mathcal{C}$. Letting $i$ be such that $\mu \in V_i$, we have that $\left|f_{\text{loc},\sigma^2}(y) - \overline{\mu}_i\right| \leq R'$. Hence, under the event in (49), by the triangle inequality,

$$|g(y) - \overline{\mu}_i| \leq R' + 8(R_0 + \sigma).$$

By assumption, $Q_0(B_{R_0}(\overline{\mu}_i)) \geq \alpha_{\min}$, so letting $R'' = R' + 8(R_0 + \sigma)$,

$$Q_{\sigma^2}(g(y) \in B_{R''}(\overline{\mu}_i)) \geq \frac{7\alpha_{\min}}{8}$$

$$Q_{\sigma^2}\left(g(y) \in \bigcup_{i=1}^{k} B_{R''}(\overline{\mu}_i)\right) \geq 1 - \frac{\alpha_{\min}}{8}.$$

By the Chernoff bounds, for independent $Z_1, \ldots, Z_N \sim \mathsf{Bernoulli}(p)$, for $c \geq 0$,

$$\mathbb{P}\left(\frac{Z_1 + \cdots + Z_N}{N} \leq (1 - c)p\right) \leq e^{-\frac{c^2 pN}{2}}$$

$$\mathbb{P}\left(\frac{Z_1 + \cdots + Z_N}{N} \geq (1 + c)p\right) \leq e^{-\frac{c^2 pN}{2+c}}.$$

For fixed $c$, it suffices to have $N = \Omega\left(\frac{\ln(1/\delta)}{p}\right)$ to make this $\leq \delta$. Applying this to $p = \frac{7\alpha_{\min}}{8}$, $(1-c)p = \frac{3\alpha_{\min}}{4}$ and $p = \frac{\alpha_{\min}}{8}$, $(1+c)p = \frac{\alpha_{\min}}{4}$ respectively, by a union bound, given $N = \Omega\left(\frac{\ln(k/\delta)}{\alpha_{\min}}\right)$ iid draws $\mu_1, \ldots, \mu_m \sim Q_{\sigma^2}$, we have

$$\mathbb{P}\left(\begin{array}{c}\forall 1 \leq i \leq k,\ |\{j : g(y_j) \in B_{R''}(\overline{\mu}_i)\}| \geq \frac{3\alpha_{\min}}{4}N \\ \text{and}\ \left|\left\{j : g(y_j) \in \bigcup_{i=1}^{k} B_{R''}(\overline{\mu}_i)\right\}\right| \geq \left(1 - \frac{\alpha_{\min}}{4}\right)N\end{array}\right) \geq 1 - \delta.$$

Suppose this event holds. Consider the sets $B_{R''+R_0}(g(y_j))$. Choose $C$ such that

$$\widetilde{R} = C\left(R_0 + \sigma\sqrt{\ln\left(\frac{1}{\alpha_{\min}}\right)}\right) \geq R'' + R_0.$$ For each $i$, choose $j(i)$ such that $g(y_{j(i)}) \in B_{R''}(\overline{\mu}_i)$. Then $B_{R_0}(\overline{\mu}_i) \subseteq B_{R''+R_0}(y_{j(i)})$ so by the second event, these $k$ sets cover $1 - \frac{\alpha_{\min}}{4}$ proportion of the $g(y_j)$, $1 \leq j \leq N$. To finish, apply Lemma 5.14 to obtain that the output of the algorithm covers $1 - \frac{\alpha_{\min}}{2}$ proportion of the $g(y_j)$. In light of the first event, for each $1 \leq i \leq k$, it must contain some $g(y_j) \in B_{R''}(\overline{\mu}_i)$. This finishes the proof. $\qquad\square$

**Lemma 5.14.** *Let $\mathcal{S}$ be a set of subsets of $X$. Let $k$ be the minimum number of sets in $\mathcal{S}$ required to cover $(1 - \varepsilon)|X|$ elements of $X$. Consider the greedy algorithm where at each step, we take the set containing the most uncovered elements, as in Algorithm 2. Then the greedy algorithm finds $O\left(k \ln\left(\frac{1}{\varepsilon}\right)\right)$ sets covering $(1 - 2\varepsilon)|X|$ elements of $X$.*

The proof is based on the classic proof of the approximation ratio for set cover [Joh73].

*Proof.* Let $S_1, S_2, \ldots$ be the sets chosen by the greedy algorithm, let $U_i = \bigcup_{i'=1}^{i} S_{i'}$, and let $x_1, x_2, \ldots$ be the covered elements in order. Define the cost of an element $x_j$ as follows: Let $S_i$ be the first set where $x_j$ appears, and set

$$c_j = \frac{1}{|S_i \setminus U_{i-1}|}.$$

Now suppose that $U_{i-1} = \{x_1, \ldots, x_{j-1}\}$, and consider the cost of $x_j$. We claim that

$$c_j \leq \frac{k}{(1 - \varepsilon)n - (j - 1)}.$$

To see this, let $A_1, \ldots, A_k$ be an optimal cover of $(1 - \varepsilon)|X|$ elements of $X$. Then they must cover $(1 - \varepsilon)|X| - (j - 1)$ elements of $X \setminus U_{i-1}$, so

$$\sum_{\ell=1}^{k} |A_\ell \cap (X \setminus U_{i-1})| \geq (1 - \varepsilon)n - (j - 1)$$

By optimality of $S_i$, $|S_i \cap (X \setminus U_{i-1})| \geq \frac{1}{k}((1 - \varepsilon)n - (j - 1))$, which shows the claim. Then the number of sets required to cover the first $(1 - 2\varepsilon)|X|$ elements is given by the sum of costs of those elements. It is the ceiling of

$$\sum_{j=1}^{\lceil(1-2\varepsilon)n\rceil} c_j \leq \sum_{j=1}^{\lceil(1-2\varepsilon)n\rceil} \frac{k}{(1 - \varepsilon)n - j + 1} \leq k \sum_{j=0}^{(1-2\varepsilon)n} \frac{1}{\varepsilon n + j} = O\left(k \ln\left(\frac{1}{\varepsilon}\right)\right). \qquad\square$$

# Acknowledgements

# References

[ABDH+18]  Hassan Ashtiani, Shai Ben-David, Nicholas Harvey, Christopher Liaw, Abbas Mehrabian, and Yaniv Plan. Nearly tight sample complexity bounds for learning mixtures of gaussians via sample compression schemes. *Advances in Neural Information Processing Systems*, 31, 2018.

[ADLS17]  Jayadev Acharya, Ilias Diakonikolas, Jerry Li, and Ludwig Schmidt. Sample-optimal density estimation in nearly-linear time. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1278–1289. SIAM, 2017.

[AK05]  Sanjeev Arora and Ravi Kannan. Learning mixtures of separated nonspherical gaussians. *Annals of Applied Probability*, pages 69–92, 2005.

[And82]  Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.

[BBPV23]  Alberto Bietti, Joan Bruna, and Loucas Pillaud-Vivien. On learning gaussian multiindex models with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023.

[BDBDD23]  Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Linear convergence bounds for diffusion models via stochastic localization. *arXiv preprint arXiv:2308.03686*, 2023.

[BDJ+22]  Ainesh Bakshi, Ilias Diakonikolas, He Jia, Daniel M Kane, Pravesh K Kothari, and Santosh S Vempala. Robustly learning mixtures of k arbitrary gaussians. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1234–1247, 2022.

[BKS99]  Itai Benjamini, Gil Kalai, and Oded Schramm. Noise sensitivity of boolean functions and applications to percolation. *Publications Mathématiques de l'Institut des Hautes Études Scientifiques*, 90:5–43, 1999.

[BM23]  Giulio Biroli and Marc Mézard. Generative diffusion in very large dimensions. *Journal of Statistical Mechanics: Theory and Experiment*, 2023(9):093402, 2023.

[CE22]  Yuansi Chen and Ronen Eldan. Localization schemes: A framework for proving mixing bounds for markov chains. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 110–122. IEEE, 2022.

[CHZW23]  Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. *arXiv preprint arXiv:2302.07194*, 2023.

[CKS24]      Sitan Chen, Vasilis Kontonis, and Kulin Shah. Learning general gaussian mixtures with efficient score matching. 2024.

[CKVEZ23]   Hugo Cui, Florent Krzakala, Eric Vanden-Eijnden, and Lenka Zdeborová. Analysis of learning a flow-based generative model from limited sample complexity. *arXiv preprint arXiv:2310.03575*, 2023.

[CL24]       Frank Cole and Yulong Lu. Score-based generative models break the curse of dimensionality in learning a family of sub-gaussian probability distributions. *arXiv preprint arXiv:2402.08082*, 2024.

[CLL23]      Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, pages 4735–4763. PMLR, 2023.

[DB22]       Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis. *arXiv preprint arXiv:2208.05314*, 2022.

[DK20]       Ilias Diakonikolas and Daniel M Kane. Small covers for near-zero sets of polynomials and learning latent variable models. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 184–195. IEEE, 2020.

[DKS17]      Ilias Diakonikolas, Daniel M Kane, and Alistair Stewart. Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 73–84. IEEE, 2017.

[DS00]       Sanjoy Dasgupta and Leonard J Schulman. A two-round variant of em for gaussian mixtures. In *Proceedings of the Sixteenth conference on Uncertainty in artificial intelligence*, pages 152–159, 2000.

[DWYZ20]    Natalie Doss, Yihong Wu, Pengkun Yang, and Harrison H Zhou. Optimal estimation of high-dimensional location gaussian mixtures. *arXiv preprint arXiv:2002.05818*, 2020.

[EAMS22]    Ahmed El Alaoui, Andrea Montanari, and Mark Sellke. Sampling from the sherrington-kirkpatrick gibbs measure via algorithmic stochastic localization. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 323–334. IEEE, 2022.

[Eld13]       Ronen Eldan. Thin shell implies spectral gap up to polylog via a stochastic localization scheme. *Geometric and Functional Analysis*, 23(2):532–569, 2013.

[GVV22]      Aparna Gupte, Neekon Vafa, and Vinod Vaikuntanathan. Continuous lwe is as hard as lwe & applications to learning gaussian mixtures. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 1162–1173. IEEE, 2022.

[HL18]        Samuel B Hopkins and Jerry Li. Mixture models, robustness, and sum of squares proofs. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1021–1034, 2018.

[IS01]     Rustam Ibragimov and Sh Sharakhmetov. The best constant in the Rosenthal inequality for nonnegative random variables. *Statistics & probability letters*, 55(4):367–376, 2001.

[Joh73]    David S Johnson. Approximation algorithms for combinatorial problems. In *Proceedings of the fifth annual ACM symposium on Theory of computing*, pages 38–49, 1973.

[KG22]     Arlene KH Kim and Adityanand Guntuboyina. Minimax bounds for estimating multivariate gaussian location mixtures. *Electronic Journal of Statistics*, 16(1):1461–1484, 2022.

[KHR22]    Frederic Koehler, Alexander Heckett, and Andrej Risteski. Statistical efficiency of score matching: The view from isoperimetry. *arXiv preprint arXiv:2210.00726*, 2022.

[KKMS08]   Adam Tauman Kalai, Adam R Klivans, Yishay Mansour, and Rocco A Servedio. Agnostically learning halfspaces. *SIAM Journal on Computing*, 37(6):1777–1805, 2008.

[KOS04]    Adam R Klivans, Ryan O'Donnell, and Rocco A Servedio. Learning intersections and thresholds of halfspaces. *Journal of Computer and System Sciences*, 68(4):808–840, 2004.

[KOS08]    Adam R Klivans, Ryan O'Donnell, and Rocco A Servedio. Learning geometric concepts via gaussian surface area. In *2008 49th Annual IEEE Symposium on Foundations of Computer Science*, pages 541–550. IEEE, 2008.

[KSS18]    Pravesh K Kothari, Jacob Steinhardt, and David Steurer. Robust moment estimation and improved clustering via sum of squares. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1035–1046, 2018.

[LC24]     Marvin Li and Sitan Chen. Critical windows: non-asymptotic theory for feature emergence in diffusion models. *arXiv preprint arXiv:2403.01633*, 2024.

[LL22]     Allen Liu and Jerry Li. Clustering mixtures with almost optimal separation in polynomial time. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1248–1261, 2022.

[LLT23]    Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general data distributions. In *International Conference on Algorithmic Learning Theory*, pages 946–985. PMLR, 2023.

[LM21]     Allen Liu and Ankur Moitra. Settling the robust learnability of mixtures of gaussians. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 518–531, 2021.

[Mon23]    Andrea Montanari. Sampling, diffusions, and stochastic localization. *arXiv preprint arXiv:2305.10690*, 2023.

[MV10]     Ankur Moitra and Gregory Valiant. Settling the polynomial learnability of mixtures of gaussians. In *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, pages 93–102. IEEE, 2010.

[MW23]     Song Mei and Yuchen Wu. Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models. *arXiv preprint arXiv:2309.11420*, 2023.

[Nel67]     Edward Nelson. *Dynamical theories of Brownian motion*, volume 101. Princeton university press, 1967.

[NP78]     SV Nagaev and IF Pinelis. Some inequalities for the distribution of sums of independent random variables. *Theory of Probability & Its Applications*, 22(2):248–256, 1978.

[OAS23]     Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution estimators. *arXiv preprint arXiv:2303.01861*, 2023.

[PRS+24]     Chirag Pabbaraju, Dhruv Rohatgi, Anish Prasad Sevekari, Holden Lee, Ankur Moitra, and Andrej Risteski. Provable benefits of score matching. *Advances in Neural Information Processing Systems*, 36, 2024.

[QR23]     Yilong Qin and Andrej Risteski. Fit like you sample: Sample-efficient generalized score matching from fast mixing markov chains. *arXiv preprint arXiv:2306.09332*, 2023.

[Rob92]     Herbert E Robbins. An empirical bayes approach to statistics. In *Breakthroughs in Statistics: Foundations and basic theory*, pages 388–394. Springer, 1992.

[RV17]     Oded Regev and Aravindan Vijayaraghavan. On learning mixtures of well-separated gaussians. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 85–96. IEEE, 2017.

[SCK23]     Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of gaussians using the ddpm objective. *arXiv preprint arXiv:2307.01178*, 2023.

[SDWMG15] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.

[SE19]     Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Proceedings of the 33rd Annual Conference on Neural Information Processing Systems*, 2019.

[SG20]     Sujayam Saha and Adityanand Guntuboyina. On the nonparametric maximum likelihood estimator for gaussian location mixture densities with application to gaussian denoising. *The Annals of Statistics*, 48(2):738–762, 2020.

[SSDK+20]     Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.

[SZ23]     Christoph Schwab and Jakob Zech. Deep learning in high dimension: Neural network expression rates for analytic functions in $L^2(\mathbb{R}^d, \gamma_d)$. *SIAM/ASA Journal on Uncertainty Quantification*, 11(1):199–234, 2023.

[Tal10]    Michel Talagrand. *Mean field models for spin glasses: Volume I: Basic examples*, volume 54. Springer Science & Business Media, 2010.

[TSM85]    David Michael Titterington, Adrian FM Smith, and Udi E Makov. Statistical analysis of finite mixture distributions. *Chichester-New York: J. Willey & Sons*, 646, 1985.

[TZ24]     Wenpin Tang and Hanyang Zhao. Score-based diffusion models via stochastic differential equations–a technical tutorial. *arXiv preprint arXiv:2402.07487*, 2024.

[Ver20]    Roman Vershynin. High-dimensional probability. *University of California, Irvine*, 10:11, 2020.

[Vin11]    Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.

[VW04]     Santosh Vempala and Grant Wang. A spectral algorithm for learning mixture models. *Journal of Computer and System Sciences*, 68(4):841–860, 2004.

[WCL+24]   Yuchen Wu, Minshuo Chen, Zihao Li, Mengdi Wang, and Yuting Wei. Theoretical insights for diffusion guidance: A case study for gaussian mixture models. *arXiv preprint arXiv:2403.01639*, 2024.

[WWY24]    Andre Wibisono, Yihong Wu, and Kaylee Yingxi Yang. Optimal score estimation via empirical bayes smoothing. *arXiv preprint arXiv:2402.07747*, 2024.

# A    Inequalities

**Theorem A.1** (Gaussian hypercontractivity, [Nel67])**.** *Let $q > p > 1$ and $t$ be such that $e^{-t} = \frac{p-1}{q-1}$. Let $g : \mathbb{R}^n \to \mathbb{R}$ be a function such that $\|g\|_{L^p(\gamma)} < \infty$. Then*

$$\|\mathscr{P}_t g\|_{L^q(\gamma)} \leq \|g\|_{L^p(\gamma)}.$$

**Lemma A.2.** *Let $q > 2$ and let $f$ be a polynomial of degree at most $d$. Then*

$$\|f\|_{L^q(\gamma)} \leq (q-1)^d \|f\|_{L^2(\gamma)}.$$

*Proof.* Write $f = \sum_{|\mathbf{k}| \leq d} a_\mathbf{k} h_\mathbf{k}$ in the Hermite basis. Then $f = \mathscr{P}_t g$ where $g = \sum_{|\mathbf{k}| \leq d} e^{|\mathbf{k}|t} a_\mathbf{k} h_\mathbf{k}$. By Gaussian hypercontractivity (Theorem A.1), choosing $t$ such that $e^{-t} = \frac{1}{q-1}$,

$$\|f\|_{L^q(\gamma)} = \|\mathscr{P}_t g\|_{L^q(\gamma)} \leq \|g\|_{L^2(\gamma)} \leq e^{td} \|f\|_{L^2(\gamma)} = (q-1)^d \|f\|_{L^2(\gamma)}. \tag{50}$$

$\square$

The following tells us how the $L^2$ norm changes when we change from one Gaussian to another, for a bounded degree polynomial.

**Lemma A.3.** *Let $f$ be a polynomial of degree at most $d$ and let $\nu$ be a measure such that for all $a \geq 0$, $\left\| \frac{d\nu}{d\gamma} \right\|_{L^{1+a}(\gamma)} \leq e^{\frac{aR^2}{2}}$ (e.g. from Lemma A.4, $Q * \mathcal{N}(0, I_n)$ where $Q$ is supported on $B_R(0)$). Then*

$$\|f\|_\nu^2 \leq \|f\|_\gamma^2 \, e^{2\sqrt{d}R}.$$

Note this works for $\mathbb{R}^n$-valued polynomials as well, since in this case $\|f\|_\nu^2 = \sum_{i=1}^n \|f_i\|_\nu^2$.

*Proof.* By Hölder's inequality and the given assumption,

$$\int_{\mathbb{R}^n} |f|^2 \, d\nu = \int_{\mathbb{R}^n} |f|^2 \frac{d\nu}{d\gamma} \, d\gamma \leq \|f\|_{L^2\left(1+\frac{1}{p}\right)(\gamma)}^2 \left\| \frac{d\nu}{d\gamma} \right\|_{L^{p+1}(\gamma)} \leq \|f\|_{L^2\left(1+\frac{1}{p}\right)(\gamma)}^2 \, e^{\frac{pR^2}{2}}.$$

By Lemma A.2 with $q = 2\left(1 + \frac{1}{p}\right)$,

$$\|f\|_{L^2\left(1+\frac{1}{p}\right)(\gamma)} \leq \left(1 + \frac{2}{p}\right)^d \|f\|_{L^2(\gamma)} \leq e^{\frac{2d}{p}} \|f\|_{L^2(\gamma)}. \tag{51}$$

Now take $p = \frac{2\sqrt{d}}{R}$ to get

$$\|f\|_\nu^2 \leq e^{\frac{2d}{p} + \frac{pR^2}{2}} \|f\|_\gamma^2 = e^{2\sqrt{d}R} \|f\|_\gamma^2.$$

$\square$

**Lemma A.4.** *Suppose $P = Q * \mathcal{N}(0, I_n)$, where $Q$ is supported on $B_R(0)$. Then for all $a \geq 0$,*

$$\left\| \frac{dP}{d\gamma} \right\|_{L^{1+a}(\gamma)} = \left( \int \left(\frac{dP}{d\gamma}\right)^{1+a} d\gamma \right)^{\frac{1}{1+a}} \leq \exp\left(\frac{aR^2}{2}\right) \quad and \quad \left\| \frac{d\gamma}{dP} \right\|_{L^{1+a}(P)} \leq \exp\left(\frac{aR^2}{2}\right)$$

*Proof.* By convexity, it suffices to consider when $P = \mathcal{N}(\mu, I_n)$ with $\|\mu\| \leq R$. Let $Z = (2\pi)^{n/2}$ be the normalizing constant of the standard Gaussian. Then

$$\left( \int \left(\frac{dP}{d\gamma}\right)^{1+a} d\gamma \right)^{\frac{1}{1+a}} = \frac{1}{Z} \int_{\mathbb{R}^n} \left( \frac{e^{-\frac{\|x\|^2}{2}}}{e^{-\frac{\|x-\mu\|^2}{2}}} \right)^a e^{-\frac{\|x-\mu\|^2}{2}} \, dx$$

$$= \frac{1}{Z} \int_{\mathbb{R}^n} e^{-\frac{\|x\|^2}{2} - (a+1)\langle x, \mu \rangle + (a+1)\frac{\|\mu\|^2}{2}}$$

$$= \frac{1}{Z} \int_{\mathbb{R}^n} e^{-\frac{\|x+(a+1)\mu\|^2}{2} + \frac{a(a+1)\|\mu\|^2}{2}} \, dx = e^{\frac{a(a+1)\|\mu\|^2}{2}} \leq e^{\frac{a(a+1)R^2}{2}}.$$

Raising to the power $\frac{1}{a+1}$ then gives the first result in this case. For $P = \mathcal{N}(\mu, I_n)$ the same inequality follows from symmetry. $\square$

# B  Sub-exponential random variables

**Definition B.1.** *We say centered random variable $X$ is $(v^2, \alpha)$ sub-exponential for $v, \alpha > 0$ if*

$$\mathbb{E}e^{\lambda X} \le e^{\frac{\lambda^2 v^2}{2}}, \forall \lambda: \quad |\lambda| \le \alpha^{-1}.$$

*We call any random variable sub-exponential if its centered version is sub-exponential.*

Sub-exponential random variables obey a desirable concentration of measure due to their fast decaying tails. This further implies their moments have a slow growth rate, a property that is of interest.

**Proposition B.2** (Sub-exponential properties [Ver20])**:** The following properties are equivalent:

- $X$ is $((c_0 v)^2, c_0 v)$-sub-exponential for some universal constant $c_0$.

- $X$ satisfies the tail bound $\mathbb{P}(X \ge t) \le 2 \exp(-c_1 t / v)$ for universal constant $c_1$.

- For all $p \ge 1$, the $p$th moment of $X$ is bounded as $\|X\|_p \le c_2 v p$ for universal constant $c_2$.

By equivalent, we mean that given one of the properties, each of the other properties holds for some constant depending on the original constant. In particular, the second property implies that the square of a $v$-sub-Gaussian variable is $(O(v^2), O(v))$-sub-exponential. Another important property of sub-exponential variables is that they behaves nicely under summation.

**Proposition B.3** (Adding sub-exponential variables)**:** Given independent $(v_i^2, \alpha_i)$-sub-exponential random variables $X_i$ for $1 \le i \le n$, $\sum_{i=1}^n X_i$ is $\left(\sum_{i=1}^n v_i^2, \max_{1 \le i \le n} \alpha_i\right)$-sub-exponential.

*Proof.* By induction, it suffices to prove this for the case $n = 2$. Our assumptions imply that $\mathbb{E}e^{\lambda X_1} \le e^{\frac{\lambda^2 v_1^2}{2}}$ and $\mathbb{E}e^{\lambda X_2} \le e^{\frac{\lambda^2 v_2^2}{2}}$ for all $\lambda$ such that $|\lambda| \le \min\{\alpha_1^{-1}, \alpha_2^{-1}\} = \max\{\alpha_1, \alpha_2\}^{-1}$. By the independence of $X_1$ and $X_2$, we therefore have

$$\mathbb{E}e^{\lambda(X_1 + X_2)} = \mathbb{E}e^{\lambda X_1}\mathbb{E}e^{\lambda X_2} \le e^{\frac{\lambda^2 v_1^2}{2}} e^{\frac{\lambda^2 v_2^2}{2}} = e^{\frac{\lambda^2\left(v_1^2 + v_2^2\right)}{2}}$$

for all such $\lambda$, so $X_1 + X_2$ is $(v_1^2 + v_2^2, \max\{\alpha_1, \alpha_2\})$, as claimed. $\square$

When the random variables are not independent, we accrue an extra factor of the number of terms.

**Proposition B.4:** Given $(v_i^2, \alpha_i)$-sub-exponential random variables $X_i$ for $1 \le i \le n$ (not necessarily independent), $\sum_{i=1}^n X_i$ is $\left(n \cdot \sum_{i=1}^n v_i^2, n \cdot \max_{1 \le i \le n} \alpha_i\right)$-sub-exponential.

*Proof.* For $|\lambda| \le \frac{1}{n} \min_{1 \le i \le n} \alpha_i^{-1} = (n \cdot \max_{1 \le i \le n} \alpha_i)^{-1}$, the defining inequality for the sub-exponential random variables are satisfied for $\lambda n$. By Hölder's inequality,

$$\mathbb{E}e^{\lambda \sum_{i=1}^n X_i} \le \prod_{i=1}^n \left(\mathbb{E}e^{\lambda n X_i}\right)^{1/n} \le \prod_{i=1}^n e^{\frac{n^2 \lambda^2 v_i^2}{2} \cdot \frac{1}{n}} = e^{\frac{\lambda^2 n \sum_{i=1}^n v_i^2}{2}}.$$

$\square$

# C Derivation of Tweedie's formula

We derive (6). Consider $\mu \sim Q$, $\xi \sim \mathcal{N}(0, \sigma^2 I_n)$ drawn independently, and let $Y = \mu + \xi$. Let $Q_{\sigma^2} = Q * \mathcal{N}(0, \sigma^2 I_n)$, and suppose it has density $q_{\sigma^2}$. By Bayes's Rule, letting $q(y|\mu) = \gamma_{\mu, \sigma^2}(y)$ denote the density of $Y$ given $\mu$,

$$
\begin{aligned}
\nabla \ln q_{\sigma^2}(y) = \nabla \log(q * \gamma_{\sigma^2}(y)) &= \nabla_y \log \int_{\mathbb{R}^d} e^{-\frac{\|y-\mu\|^2}{2\sigma^2}} \, dQ_0(\mu) \\
&= -\frac{\int_{\mathbb{R}^d} \frac{y-\mu}{\sigma^2} e^{-\frac{\|y-\mu\|^2}{2\sigma^2}} \, dQ_0(\mu)}{\int_{\mathbb{R}^d} e^{-\frac{\|y-\mu\|^2}{2\sigma^2}} \, dQ_0(\mu)} = \frac{\int_{\mathbb{R}^d} \frac{\mu-y}{\sigma^2} q(y|\mu) \, dQ_0(\mu)}{\int_{\mathbb{R}^d} q(y|\mu) \, dQ_0(\mu)} \\
&= \frac{1}{\sigma^2} \mathbb{E}[\mu - y | Y = y].
\end{aligned}
$$