

Efficient Truncated Statistics with Unknown Truncation

Vasilis Kontonis
University of Wisconsin-Madison
 Madison, WI, USA
 kontonis@wisc.edu

Christos Tzamos
University of Wisconsin-Madison
 Madison, WI, USA
 tzamos@wisc.edu

Manolis Zampetakis
 MIT
 Cambridge, MA, USA
 mzampet@mit.edu

Abstract—We study the problem of estimating the parameters of a Gaussian distribution when samples are only shown if they fall in some (unknown) set. This core problem in truncated statistics has long history going back to Galton, Lee, Pearson and Fisher. Recent work by Daskalakis et al. (FOCS’18), provides the first efficient algorithm that works for arbitrary sets in high dimension when the set is known, but leaves as an open problem the more challenging and relevant case of unknown truncation set.

Our main result is a computationally and sample efficient algorithm for estimating the parameters of the Gaussian under arbitrary unknown truncation sets whose performance decays with a natural measure of complexity of the set, namely its Gaussian surface area. Notably, this algorithm works for large families of sets including intersections of halfspaces, polynomial threshold functions and general convex sets. We show that our algorithm closely captures the tradeoff between the complexity of the set and the number of samples needed to learn the parameters by exhibiting a set with small Gaussian surface area for which it is information theoretically impossible to learn the true Gaussian with few samples.

Keywords-Truncated Statistics; Gaussian Distribution; Learning Theory; Unknown Truncation Set

I. INTRODUCTION

A classical challenge in Statistics is estimation from truncated samples. Truncation occurs when samples falling outside of some subset S of the support of the distribution are not observed. Truncation of samples has myriad manifestations in business, economics, engineering, social sciences, and all areas of the physical sciences.

Statistical estimation under truncated samples has had a long history in Statistics, going back to at least the work of Galton [Gal97] who analyzed truncated samples corresponding to speeds of American trotting horses. Following Galton’s work, Pearson and Lee [Pea02], [PL08], [Lee14] used the method of moments in order to estimate the mean and standard deviation of a truncated univariate normal distribution and later Fisher [Fis31] used the maximum likelihood method for the same estimation problem. Since then, there has been a large volume of research devoted to estimating the truncated normal distribution; see e.g. [Sch86], [Coh16], [BC14]. Nevertheless, the first algorithm that is provably computationally and statistically efficient was only recently developed by Daskalakis et al. [DGTZ18], under the assumption that the truncation set S is known.

In virtually all these works the question of estimation under unknown truncation set is raised. Our work resolves this question by providing tight sample complexity guarantees and an efficient algorithm for recovering the underlying Gaussian distribution. Although this estimation problem has clear and important practical and theoretical motivation too little was known prior to our work even in the asymptotic regime. In the early work of Shah and Jaiswal [SJ66] it was proven that the method of moments can be used to estimate a single dimensional Gaussian distribution when the truncation set is unknown but it is assumed to be an interval. In the other extreme where the set is allowed to be arbitrarily complex, Daskalakis et al. [DGTZ18] showed that it is information theoretically impossible to recover the parameters. We provide the first complete analysis of the number of samples needed for recovery taking into account the complexity of the underlying set.

Our Contributions.: Our work studies the estimation task when the truncation set belongs in a family \mathcal{C} of “low complexity”. We use two different notions for quantifying the complexity of sets: the VC-dimension and the Gaussian Surface Area.

Our first result is that for any set family with VC-dimension $\text{VC}(\mathcal{C})$, the mean and covariance of the true d -dimensional Gaussian Distribution can be recovered up to accuracy ε using only $\tilde{O}\left(\frac{\text{VC}(\mathcal{C})}{\varepsilon} + \frac{d^2}{\varepsilon^2}\right)$ truncated samples.

Informal Theorem 1. *Let \mathcal{C} be a class of sets with VC-dimension $\text{VC}(\mathcal{C})$ and let $N = \tilde{O}\left(\frac{\text{VC}(\mathcal{C})}{\varepsilon} + \frac{d^2}{\varepsilon^2}\right)$. Given N samples from a d -dimensional Gaussian $\mathcal{N}(\mu, \Sigma)$ with unknown mean μ and covariance Σ , truncated on a set $S \in \mathcal{C}$ with mass at least α , it is possible to find an estimate $(\hat{\mu}, \hat{\Sigma})$ such that $d_{\text{TV}}(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\hat{\mu}, \hat{\Sigma})) \leq \varepsilon$.*

The estimation method computes the set of smallest mass that maximizes the likelihood of the data observed and learns the truncated distribution within error $O(\varepsilon)$ in total variation distance. To translate this error in total variation to parameter distance, we prove a general result showing that it is impossible to create a set (no matter the complexity) so that two Gaussians whose parameters are far have similar truncated distributions (see Lemma 8).

A simple but not successful approach would be to first

try to learn an approximation of the truncation set with symmetric difference roughly ε^2/d^2 with the true set and then run the algorithm of [DGTZ18] using the approximate oracle. This approach would lead to a $\text{VC}(\mathcal{S})d^2/\varepsilon^2$ sample complexity that is worse than what we get. More importantly, doing empirical risk minimization¹ using truncated samples does not guarantee that we will find a set of small symmetric difference with the true and it is not clear how one could achieve that.

Our result bounds the sample complexity of identifying the underlying Gaussian distribution in terms of the VC-dimension of the set but does not yield a computationally efficient method for recovery. Obtaining a computationally efficient algorithm seems unlikely, unless one restricts attention to simple specific set families, such as axis aligned rectangles. One would hope that exploiting the fact that samples are drawn from a “tame” distribution, such as a Gaussian, can lead to general computationally efficient algorithms and even improved sample complexity.

Indeed, our main result is an algorithm that is both computationally and statistically efficient for estimating the parameters of a spherical Gaussian and uses only $d^{O(\Gamma(\mathcal{C}))}$ samples, where $\Gamma(\mathcal{C})$ is the *Gaussian Surface Area* of the class \mathcal{C} , an alternative complexity measure introduced by Klivans et al. [KOS08]:

Informal Theorem 2. *Let \mathcal{C} be a class of sets with Gaussian surface area at most $\Gamma(\mathcal{C})$ and let $k = \text{poly}(1/\alpha, 1/\varepsilon)\Gamma(\mathcal{C})^2$. Given $N = d^k$ samples from a spherical d -dimensional Gaussian $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$, truncated on a set $S \in \mathcal{C}$ with mass at least α , in time $\text{poly}(m)$, we can find an estimate $\hat{\boldsymbol{\mu}}, \hat{\sigma}^2$ such that*

$$d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I}), \mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\sigma}^2 \mathbf{I})) \leq \varepsilon.$$

The notion of Gaussian surface area can lead to better sample complexity bounds even when the VC dimension is infinite. An example of such a case is when \mathcal{C} is the class of all convex sets. Table I summarizes the known bounds for the Gaussian surface area of different concept classes and the implied sample complexity in our setting when combined with our main theorem.

Beyond spherical Gaussians, our main result extends to Gaussians with arbitrary diagonal covariance matrices. In addition, we provide an information theoretic result showing that the case with general covariance matrices can also be estimated using the same sample complexity bound by finding a Gaussian and a set that matches the moments of the true distribution. We remark our main algorithmic result Informal Theorem 3 uses Gaussian Surface Area whereas our sample complexity result Informal Theorem 2 uses VC-dimension. We discuss the differences of the two approaches in Section VII.

¹That is finding a set of the family that contains all the observed samples.

Informal Theorem 3. *Let \mathcal{C} be a class of sets with Gaussian surface area at most $\Gamma(\mathcal{C})$ and let $k = \text{poly}(1/\alpha, 1/\varepsilon)\Gamma(\mathcal{C})^2$. Any truncated Gaussian with $\mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}, \hat{S})$ with $\hat{S} \in \mathcal{C}$ that approximately matches the moments up to degree k of a truncated d -dimensional Gaussian $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ with $S \in \mathcal{C}$, satisfies $d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})) \leq \varepsilon$. The number of samples to estimate the moments within the required accuracy is at most $d^{O(k)}$.*

This shows that the first few moments are sufficient to identify the parameters. Analyzing the guarantees of moment matching methods is notoriously challenging as it involves bounding the error of a system of many polynomial equations. Even for a single-dimensional Gaussian with truncation in an interval, where closed form solutions of the moments exist, it is highly non-trivial to bound these errors [SJ66]. In contrast, our analysis using Hermite polynomials allows us to easily obtain bounds for arbitrary truncation sets in high dimensions, even though no closed form expression for the moments exists.

We conclude by showing that the dependence of our sample complexity bounds both on the VC-dimension and the Gaussian Surface Area is tight up to polynomial factors. In particular, we construct a family in d dimensions with VC dimension 2^d and Gaussian surface area $O(d)$ for which it is not possible to learn the mean of the underlying Gaussian within 1 standard deviation using $o(2^{d/2})$ samples.

Informal Theorem 4. *There exists a family of sets \mathcal{S} with $\Gamma(\mathcal{S}) = O(d)$ and VC-dimension 2^d such that any algorithm that draws N samples from $\mathcal{N}(\boldsymbol{\mu}, \mathbf{I}, S)$ and computes an estimate $\tilde{\boldsymbol{\mu}}$ with $\|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2 \leq 1$ must have $N = \Omega(2^{d/2})$.*

Our techniques and relation to prior work.: The work of Klivans et al. [KOS08] provides a computationally and sample efficient algorithm for learning geometric concepts from labeled examples drawn from a Gaussian distribution. On the other hand, the recent work of Daskalakis et al. [DGTZ18] provides efficient estimators for truncated statistics with *known* sets. One could hope to combine these two approaches for our setting, by first learning the set and then using the algorithm of [DGTZ18] to learn the parameters of the Gaussian. This approach, however, fails for two reasons. First, the results of Klivans et al. [KOS08] apply in the supervised learning setting where one has access to both positive and negative samples, while our problem can be thought of as observing only positive examples (those falling inside the set). In addition, any direct approach that extends their result to work with positive only examples requires that the underlying Gaussian distribution is known in advance.

One of our key technical contributions is to extend the techniques of Klivans et al. [KOS08] to work with *positive only examples* from an *unknown* Gaussian distribution, which is the major case of interest in truncated statistics. To

Concept Class	Gaussian Surface Area	Sample Complexity
Polynomial threshold functions of degree k	$O(k)$ [Kan11]	$d^{O(k^2)}$
Intersections of k halfspaces	$O(\sqrt{\log k})$ [KOS08]	$d^{O(\log k)}$
General convex sets	$O(d^{1/4})$ [Bal93]	$d^{O(\sqrt{d})}$

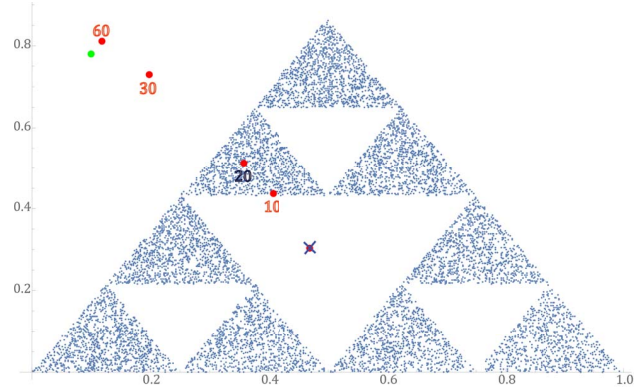
Table I
SUMMARY OF KNOWN RESULTS FOR GAUSSIAN SURFACE AREA. THE LAST COLUMN GIVES THE SAMPLE COMPLEXITY WE OBTAIN FOR OUR SETTING.

perform the set estimation Klivans et al. [KOS08], rely on a family of orthogonal polynomials with respect to the Gaussian distribution, namely the Hermite polynomials and show that the indicator function of the set is well approximated by its low degree Hermite expansion. While we cannot learn this function directly in our setting, we are able to recover an alternative function, that contains “entangled” information of both the true Gaussian parameters and the underlying set. After learning the function, we formulate an optimization problem whose solution enables us to decouple these two quantities and retrieve both the Gaussian parameters and the underlying set. We describe our estimation method in more detail in Section IV. As a corollary of our approach, we obtain the first efficient algorithm for learning geometric concepts from positive examples drawn from an unknown spherical Gaussian distribution.

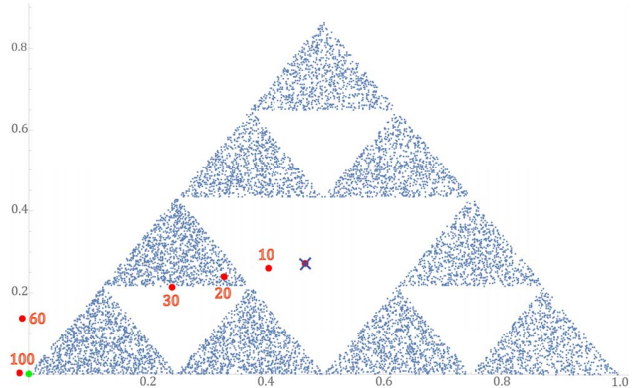
Simulations.: In addition to the theoretical guarantees of our algorithm, we empirically evaluate its performance using simulated data. We present the results that we get in Figure 1, where one can see that even when the truncation set is complex, our algorithm finds an accurate estimation of the mean of the untruncated distribution. Observe that our algorithm succeeds in estimating the true mean of the input distribution despite the fact that the set is unknown and the samples look similar in both cases.

A. Further Related Work

Our work is related to the field of robust statistics as it can robustly learn a Gaussian even in the presence of an adversary erasing samples outside a certain set. Recently, there has been a lot of theoretical work doing robust estimation of the parameters of multi-variate Gaussian distributions in the presence of arbitrary corruptions to a small ε fraction of the samples, allowing for both deletions of samples and additions of samples that can also be chosen adaptively [DKK⁺16], [CSV17], [LRV16], [DKK⁺17], [DKK⁺18]. When the corruption of the data is so powerful it is easy to see that the estimation error of the parameter depends on ε and cannot shrink to 0 as the number of samples grows to infinity. In our model the corruption is more restrictive but in return our results show how to estimate the parameters of a multi-variate Gaussian distribution to arbitrary accuracy even when the fraction of corruption is any constant less than 1.



(a) Execution of our algorithm for isotropic Gaussian distribution with $\mu^* = (0.1, 0.78)$ and $\mu_S = (0.48, 0.32)$.



(b) Execution of our algorithm for isotropic Gaussian distribution with $\mu^* = (0, 0)$ and $\mu_S = (0.47, 0.27)$.

Figure 1. Illustration of the results of our algorithm for an unknown truncation set. The \times sign corresponds to the conditional mean of the truncated distribution, while the green point corresponds to the true mean and the red points correspond to the estimated true mean depending on the degree of the Hermite polynomials that are being used by the algorithm.

Our work also has connections with the literature of learning from positive examples. At the heart of virtually all of the results in this literature is the use of the exact knowledge of the original non-truncated distribution to be able to generate fake negative examples, e.g. [Den98], [LDG00]. When the original distribution is uniform, better algorithms are known. Diakonikolas et al. [DDS14] gave

efficient learning algorithms for DNFs and linear threshold functions, Frieze et al. [FJK96] and Anderson et al. [AGR13] gave efficient learning algorithms for learning d -dimensional simplices. Another line of work proves lower bounds on the sample complexity of recovering an unknown set from positive examples. Goyal et al. [GR09] showed that learning a convex set in d -dimensions to accuracy ε from positive samples, uniformly distributed inside the set, requires at least $2^{\Omega(\sqrt{d}/\varepsilon)}$ samples, while the work of [Eld11] showed that $2^{\Omega(\sqrt{d})}$ samples are necessary even to estimate the mass of the set. To the best of our knowledge, no matching upper bounds are known for those results. Our estimation result implies that $d^{\text{poly}(\frac{1}{\varepsilon})\sqrt{d}}$ are sufficient to learn the set and its mass when given positive samples from a Gaussian truncated on the convex set.

II. PRELIMINARIES

Notation. We use small bold letters \mathbf{x} to refer to real vectors in finite dimension \mathbb{R}^d and capital bold letters \mathbf{A} to refer to matrices in $\mathbb{R}^{d \times \ell}$. Similarly, a function with image in \mathbb{R}^d is represented by a small and bold letter \mathbf{f} . Given a subset S of \mathbb{R}^d we define $\mathbf{1}_S(\mathbf{x})$ to be its 0 – 1 indicator. Let $\mathbf{A} \in \mathbb{R}^{d \times d}$, we define $\mathbf{A}^b \in \mathbb{R}^{d^2}$ to be the standard vectorization of \mathbf{A} . Let also \mathcal{Q}_d be the set of all the symmetric $d \times d$ matrices. The *Frobenius norm* of a matrix \mathbf{A} is defined as $\|\mathbf{A}\|_F = \|\mathbf{A}^b\|_2$.

Gaussian Distribution. Let $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be the normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, with the following probability density function

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{x}) = \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (1)$$

Also, let $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; S)$ denote the *probability mass* of a measurable set S under this Gaussian measure. We shall also denote by \mathcal{N}_0 the standard Gaussian, whether it is single or multidimensional will be clear from the context.

Truncated Gaussian Distribution. Let $S \subseteq \mathbb{R}^d$ be a subset of the d -dimensional Euclidean space, we define the *S -truncated normal distribution* $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ the normal distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ conditioned on taking values in the subset S . The probability density function of $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ is the following

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S; \mathbf{x}) = \frac{\mathbf{1}_S(\mathbf{x})}{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; S)} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{x}). \quad (2)$$

We will assume that the covariance matrix $\boldsymbol{\Sigma}$ is full rank. The case where $\boldsymbol{\Sigma}$ is not full rank we can easily detect and solve the estimation problem in the linear subspace of samples.

The core complexity measure of Borel sets in \mathbb{R}^d that we use is the notion of Gaussian Surface Area defined below.

Definition 1 (GAUSSIAN SURFACE AREA). For a Borel set $A \subseteq \mathbb{R}^d$, $\delta \geq 0$ let $A_\delta = \{x : \text{dist}(x, A) \leq \delta\}$. The *Gaussian surface area* of A is

$$\Gamma(A) = \liminf_{\delta \rightarrow 0} \frac{\mathcal{N}_0(A_\delta \setminus A)}{\delta}.$$

We define the *Gaussian surface area* of a family of sets \mathcal{C} to be $\Gamma(\mathcal{C}) = \sup_{C \in \mathcal{C}} \Gamma(C)$.

Problem Formulation. Given samples from a truncated Gaussian $\mathcal{N}_S^* := \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$, our goal is to learn the parameters $(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$ and recover the set S . We denote by $\alpha^* = \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; S)$, the total mass contained in set S by the untruncated Gaussian $\mathcal{N}^* := \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$. Throughout this paper, we assume that we know an absolute constant $\alpha > 0$ such that

$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; S) = \alpha^* \geq \alpha. \quad (3)$$

We next state the following simple lemma that connects the total variation distance of two Normal distributions with their parameter distance. For a proof see e.g. Corollaries 2.13 and 2.14 of [DKK⁺16].

Lemma 1. Let $N_1 = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$, $N_2 = \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ be two Normal distributions. Then

$$d_{\text{TV}}(N_1, N_2) \leq \frac{1}{2} \left\| \boldsymbol{\Sigma}_1^{-1/2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right\|_2 + \sqrt{2} \left\| \mathbf{I} - \boldsymbol{\Sigma}_1^{-1/2} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_1^{-1/2} \right\|_F$$

We readily use the following two lemmas from [DGTZ18]. The first suggests that we can accurately estimate the parameters $(\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S)$.

Lemma 2. Let $(\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S)$ be the mean and covariance of the truncated Gaussian $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ for a set S such that $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; S) \geq \alpha$. Using $\tilde{O}(\frac{d}{\varepsilon^2} \log(1/\alpha) \log^2(1/\delta))$ samples, we can compute estimates $\tilde{\boldsymbol{\mu}}_S$ and $\tilde{\boldsymbol{\Sigma}}_S$ such that, with probability at least $1 - \delta$,

$$\begin{aligned} \|\boldsymbol{\Sigma}^{-1/2}(\tilde{\boldsymbol{\mu}}_S - \boldsymbol{\mu}_S)\|_2 &\leq \varepsilon \\ (1 - \varepsilon)\boldsymbol{\Sigma}_S &\preceq \tilde{\boldsymbol{\Sigma}}_S \preceq (1 + \varepsilon)\boldsymbol{\Sigma}_S \end{aligned}$$

The second lemma suggests that the empirical estimates are close to the true parameters of underlying truncated Gaussian.

Lemma 3. The empirical mean and covariance $\tilde{\boldsymbol{\mu}}_S$ and $\tilde{\boldsymbol{\Sigma}}_S$ computed using $\tilde{O}(d^2 \log^2(1/\alpha\delta))$ samples from a truncated Normal $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ with $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; S) \geq \alpha$ satisfies with probability $1 - \delta$ that:

$$\begin{aligned} \|\boldsymbol{\Sigma}^{-1/2}(\tilde{\boldsymbol{\mu}}_S - \boldsymbol{\mu})\|_2^2 &\leq O(\log \frac{1}{\alpha}), \quad \tilde{\boldsymbol{\Sigma}}_S \succeq \Omega(\alpha^2)\boldsymbol{\Sigma}, \\ \left\| \boldsymbol{\Sigma}^{-1/2} \tilde{\boldsymbol{\Sigma}}_S \boldsymbol{\Sigma}^{-1/2} - \mathbf{I} \right\|_F^2 &\leq O(\log \frac{1}{\alpha}). \end{aligned}$$

Moreover, $\Omega(\alpha^2) \leq \left\| \tilde{\boldsymbol{\Sigma}}_S^{-1/2} \boldsymbol{\Sigma} \tilde{\boldsymbol{\Sigma}}_S^{-1/2} \right\|_2 \leq O(1/\alpha^2)$.

In particular, the mean and covariance $\tilde{\boldsymbol{\mu}}_S$ and $\tilde{\boldsymbol{\Sigma}}_S$ that satisfy the conditions of Lemma 3, are in $(O(\log(1/\alpha)), 1 - O(\alpha^2))$ -near isotropic position.

We will use the following very useful anti-concentration result about the Gaussian mass of sets defined by polynomials.

Theorem 1 (Theorem 8 of [CW01]). *Let $q, \gamma \in \mathbb{R}_+$, $\boldsymbol{\mu} \in \mathbb{R}^d$, $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ such that $\boldsymbol{\Sigma}$ is symmetric positive semidefinite and $p : \mathbb{R}^d \rightarrow \mathbb{R}$ be a multivariate polynomial of degree at most ℓ , we define*

$$\bar{Q} = \{\mathbf{x} \in \mathbb{R}^d \mid |p(\mathbf{x})| \leq \gamma\},$$

then there exists an absolute constant C such that

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \bar{Q}) \leq \frac{Cq\gamma^{1/\ell}}{\left(\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [|p(\mathbf{z})|^{q/\ell}]\right)^{1/q}}.$$

A. Hermite Polynomials, Ornstein-Uhlenbeck Operator, and Gaussian Surface Area.

We denote by $L^2(\mathbb{R}^d, \mathcal{N}_0)$ the vector space of all functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0}[f^2(\mathbf{x})] < \infty$. The usual inner product for this space is $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0}[f(\mathbf{x})g(\mathbf{x})]$. While, usually one considers the probabilists' or physicists' Hermite polynomials, in this work we define the *normalized* Hermite polynomial of degree i to be $H_0(x) = 1, H_1(x) = x, H_2(x) = \frac{x^2-1}{\sqrt{2}}, \dots, H_i(x) = \frac{He_i(x)}{\sqrt{i!}}, \dots$ where by $He_i(x)$ we denote the probabilists' Hermite polynomial of degree i . These normalized Hermite polynomials form a complete orthonormal basis for the single dimensional version of the inner product space defined above. To get an orthonormal basis for $L^2(\mathbb{R}^d, \mathcal{N}_0)$, we use a multi-index $V \in \mathbb{N}^d$ to define the d -variate normalized Hermite polynomial as $H_V(\mathbf{x}) = \prod_{i=1}^d H_{v_i}(x_i)$. The total degree of H_V is $|V| = \sum v_i \in V v_i$. Given a function $f \in L^2$ we compute its Hermite coefficients as $\hat{f}(V) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0}[f(\mathbf{x})H_V(\mathbf{x})]$ and express it uniquely as $\sum_{V \in \mathbb{N}^d} \hat{f}(V)H_V(\mathbf{x})$. We denote by $S_k f(x)$ the degree k partial sum of the Hermite expansion of f , $S_k f(\mathbf{x}) = \sum_{|V| \leq k} \hat{f}(V)H_V(\mathbf{x})$. Then, since the basis of Hermite polynomials is complete, we have $\lim_{k \rightarrow \infty} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0}[(f(\mathbf{x}) - S_k f(\mathbf{x}))^2] = 0$. We would like to quantify the convergence rate of $S_k f$ to f . Parseval's identity states that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(f(\mathbf{x}) - S_k f(\mathbf{x}))^2] = \sum_{|V|=k}^{\infty} \hat{f}(V)^2.$$

Definition 2 (HERMITE CONCENTRATION). *Let $\gamma(\varepsilon, d)$ be a function $\gamma : (0, 1/2) \times \mathbb{N} \mapsto \mathbb{N}$. We say that a class of functions \mathcal{F} over \mathbb{R}^d has a Hermite concentration bound of $\gamma(\varepsilon, d)$, if for all $d \geq 1$, all $\varepsilon \in (0, 1/2)$, and $f \in \mathcal{F}$ it holds $\sum_{|V| \geq \gamma(\varepsilon, d)} \hat{f}(V)^2 \leq \varepsilon$.*

We now define the Gaussian Noise Operator as in [O'D14]. Using a different parametrization, which is not

convenient for our purposes, these operators are also known as the Ornstein-Uhlenbeck semigroup, or the Mehler transform.

Definition 3. *The Gaussian Noise operator T_ρ is the linear operator defined on the space of functions $L^1(\mathbb{R}^d, \mathcal{N}_0)$ by*

$$T_\rho f(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim \mathcal{N}_0} [f(\rho\mathbf{x} + \sqrt{1 - \rho^2}\mathbf{y})].$$

A nice property of operator $T_{1-\rho}$ that we will use is that it has a simple Hermite expansion

$$S_k(T_\rho f)(\mathbf{x}) = \sum_{V: |V| \leq k} \rho^{|V|} \hat{f}(V) H_V(\mathbf{x}) \quad (4)$$

We also define the noise sensitivity of a function f .

Definition 4 (NOISE SENSITIVITY). *Let $f : \mathbb{R}^d \mapsto \mathbb{R}$ be a function in $L^2(\mathbb{R}^d, \mathcal{N}_0)$. The noise sensitivity of f at $\rho \in [0, 1]$ is defined to be*

$$\mathbf{NS}_\rho[f] = 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [f(\mathbf{x})^2 - f(\mathbf{x})T_{1-\rho}f(\mathbf{x})]$$

Since, the vectors \mathbf{x} and $\mathbf{z} = (1 - \rho)\mathbf{x} + \sqrt{1 - \rho^2}\mathbf{y}$ are jointly distributed according to

$$D_\rho = \mathcal{N} \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{I} & (1-\rho)\mathbf{I} \\ (1-\rho)\mathbf{I} & \mathbf{I} \end{pmatrix} \right). \quad (5)$$

we can write

$$\mathbf{NS}_\rho[f] \quad (6)$$

$$= \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [f(\mathbf{x})^2] + \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [f(\mathbf{z})^2 - 2f(\mathbf{x})f(\mathbf{z})] \\ = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [(f(\mathbf{x}) - f(\mathbf{z}))^2]. \quad (7)$$

When f is an indicator of a set, the noise sensitivity is

$$\mathbf{NS}_\rho[\mathbf{1}_S] = 2 \mathbb{E}_{(\mathbf{x}, \mathbf{z})} [\mathbf{1}_S(\mathbf{x})(1 - \mathbf{1}_S(\mathbf{z}))] \\ = 2 \mathbb{E}_{(\mathbf{x}, \mathbf{z})} [\mathbf{1}_S(\mathbf{x})\mathbf{1}_{S^c}(\mathbf{z})], \quad (8)$$

which is equal to the probability of the correlated points \mathbf{x}, \mathbf{z} landing at "opposite" sides of S .

Ledoux [Led94] and Pisier [Pis86] showed that the noise sensitivity of a set can be bounded by its Gaussian surface area.

Definition 5 (Gaussian Surface Area). *For a Borel set $A \subseteq \mathbb{R}^d$, its Gaussian surface area is $\Gamma(A) = \liminf_{\delta \rightarrow 0} \frac{\mathcal{N}_0(A_\delta \setminus A)}{\delta}$, where $A_\delta = \{x : \text{dist}(x, A) \leq \delta\}$.*

We will use the following lemma given in [KOS08].

Lemma 4 (Corollary 14 of [KOS08]). *For a Borel set $S \subseteq \mathbb{R}^d$ and $\rho \geq 0$, $\mathbf{NS}_\rho[\mathbf{1}_S(\mathbf{x})] \leq \sqrt{\pi} \sqrt{\rho} \Gamma(S)$.*

For more details on the Gaussian space and Hermite Analysis (especially from the theoretical computer science perspective), we refer the reader to [O'D14]. Most of the facts about Hermite polynomials that we shall use in this work are well known properties and can be found, for example, in [Sze67].

III. IDENTIFIABILITY WITH BOUNDED VC DIMENSION

In this section we analyze the sample complexity of learning the true Gaussian parameters when the truncation set has bounded VC-dimension. In particular we show that the overhead over the d^2/ε^2 samples (which is the sample complexity of learning the parameters of the Gaussian without truncation) is proportional to the VC dimension of the class.

Theorem 2. *Let \mathcal{S} be a family of sets of finite VC dimension, and let $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ be a truncated Gaussian distribution such that $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; S) \geq \alpha$. Given N samples with*

$$N = \text{poly}(1/\alpha) \tilde{O}\left(\frac{d^2}{\varepsilon^2} + \frac{\text{VC}(\mathcal{S})}{\varepsilon}\right)$$

Then, with probability at least 99%, it is possible to identify $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ that satisfy $d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})) \leq \varepsilon$ and $\|\boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}})\|_2 \leq \varepsilon$ and $\|\mathbf{I} - \boldsymbol{\Sigma}^{-1/2}\tilde{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1/2}\|_F \leq \varepsilon$.

Our algorithm works by first learning the truncated distribution within total variation distance ε . To do this, we first assume that we know the mean and covariance of the underlying Gaussian by guessing the parameters and accurately learn the underlying set. After drawing $N = \Theta\left(\frac{\text{VC}(\mathcal{S})\log(1/\varepsilon)}{\varepsilon}\right)$ samples from the distribution, any set in the class that contains the samples will only exclude at most an ε fraction of the total mass. Picking the set \tilde{S} that maximizes the likelihood of those samples, i.e. the set with minimum mass according to the guessed Gaussian distribution, guarantees that the total variation distance between the learned truncated distribution and the true is at most ε , if the guess of the parameters was accurate (Lemma 5). For the proof of Lemma 5 we refer the reader to the full version of the paper.

Lemma 5. *Let \mathcal{S} be a family of subsets in \mathbb{R}^d and Let $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S^*) = \mathcal{N}_{\tilde{S}}^*$ be a Normal distribution truncated on the set $S^* \in \mathcal{S}$. Fix $\varepsilon \in (0, 1)$, $\delta \in (0, 1/4)$ and let*

$$N = O\left(\frac{\text{VC}(\mathcal{S})\log(1/\varepsilon)}{\varepsilon} + \log\left(\frac{1}{\delta}\right)\right)$$

Moreover, let $\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}$ be such that $d_{\text{TV}}(\mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}), \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})) \leq \varepsilon$. Assume that we draw N samples \mathbf{x}_i from \mathcal{N}_{S^} , Let \tilde{S} be the solution of the problem*

$$\min_S \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}; S) \text{ subject to } \mathbf{x}_i \in S \text{ for all } i \in [n]$$

Then with probability at least $1 - \delta$ we have $d_{\text{TV}}(\mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}, \tilde{S}), \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)) \leq 3\varepsilon/(2\alpha)$.

This is because the total variation distance between two densities f and g can be written as $\int (f(x) - g(x))\mathbf{1}_{f(x) > g(x)} dx$. Note that by choosing the set of the smallest mass consistent with the samples, we guarantee that the guess will have higher density at every point apart from those outside the support \tilde{S} . However, as we argued the

outside mass is at most ε with respect to the true distribution which gives the bound in the total variation distance.

To remove the assumption that the true parameters are known, we build a cover of all possible mean and covariance matrices that the underlying Gaussian might have and run the tournament from [DK14] to identify the best one (Lemma 6). We will use the following statement of the tournament from [DK14]. See also [DL12].

Lemma 6 (Tournament [DK14]). *There is an algorithm, which is given sample access to some distribution X and a collection of distributions $\mathcal{H} = \{H_1, \dots, H_N\}$ over some set, access to a PDF comparator for every pair of distributions $H_i, H_j \in \mathcal{H}$, an accuracy parameter $\varepsilon > 0$, and a confidence parameter $\delta > 0$. The algorithm makes $O(\log(1/\delta)\varepsilon^2 \log N)$ draws from each of X, H_1, \dots, H_N and returns some $H \in \mathcal{H}$ or declares "failure" If there is some $H \in \mathcal{H}$ such that $d_{\text{TV}}(H, X) \leq \varepsilon$ then with probability at least $1 - \delta$ the returned distribution H satisfies $d_{\text{TV}}(H, X) \leq 512\varepsilon$. The total number of operations of the algorithm is $O(\log(1/\delta)(1/\varepsilon^2)(N \log N + \log 1/\delta))$.*

While there are $(d/\varepsilon)^{O(d^2)}$ such parameters, the number of samples needed for running the tournament is only logarithmic which shows that an additional $\tilde{O}(d^2/\varepsilon^2)$ are sufficient to find a hypothesis in total variation distance ε (Lemma 7). The proof of Lemma 7 can be found the full version of the paper.

Lemma 7. *Let $S \in \mathcal{S}$ be a subset of \mathbb{R}^d and $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ be the corresponding truncated normal distribution. Then $\tilde{O}(\text{VC}(\mathcal{S})/\varepsilon + d^2/\varepsilon^2)$ samples are sufficient to find parameters $\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}, \tilde{S}$ such that $d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S), \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}, \tilde{S})) \leq \varepsilon$ with probability at least 99%.*

We finally argue that the ε error in total variation of the truncated distributions translates to an $O(\varepsilon)$ bound in total variation distance of the untruncated distributions (Lemma 8). We show that this is true in general and does not depend on the complexity of the set. To prove this statement, we consider two Gaussians with parameters that are far from each other and construct the worst possible set to make their truncated distributions as close as possible. We show that under the requirement that the set contains at least α mass, the total variation distance of the truncated distributions will be large.

Lemma 8 (Total Variation of Truncated Normals). *Let $D_1 = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, S_1)$ and $D_2 = \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, S_2)$ be two truncated Normal distributions such that $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1; S_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2; S_2) \geq \alpha$. Then*

$$d_{\text{TV}}(D_1, D_2) \geq C_\alpha d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2))$$

where $C_\alpha < \alpha/8$ is a positive constant that only depends on α , $C_\alpha = \Omega(\alpha^3)$.

Proof: Without loss of generality we assume that $D_1 = \mathcal{N}(\mathbf{0}, \mathbf{I}, S_1)$ and $D_2 = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}, S_2)$, where $\boldsymbol{\Lambda}$ is a diagonal matrix. We want to find the worst sets S_1, S_2 so that $d_{\text{TV}}(D_1, D_2)$ is small. If $D_1(S_1 \setminus S_2) \geq \alpha/2$ then the statement holds. Therefore, we consider the set $S = S_1 \cap S_2$ and relax the constraint that the truncated Gaussian D_2 integrates to 1. Taking into account the fact that the set $S = S_1 \cap S_2$ must have at least some mass $\alpha/2$ with respect to $\mathcal{N}(\mathbf{0}, \mathbf{I})$, the following optimization problem provides a lower bound on the total variation distance of D_1 and D_2 .

$$\begin{aligned} \min_{S \in \mathcal{S}, \beta > 0} \quad & \frac{1}{\alpha} \int |\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) - \frac{\alpha}{\beta} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}; \mathbf{x})| \mathbf{1}_S(\mathbf{x}) d\mathbf{x} \\ \text{subj. to} \quad & \int \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) \mathbf{1}_S(\mathbf{x}) d\mathbf{x} \geq \alpha/2, \end{aligned}$$

For any fixed $\beta > 0$ this is a fractional knapsack problem and therefore we should include in the set the points \mathbf{x} in order of increasing ratio of weight that is contribution to the L_1 error $|\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) - \frac{\alpha}{\beta} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}; \mathbf{x})|$, over value, that is density $\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})$ until we reach some threshold T . Therefore, the set is defined to be

$$\begin{aligned} S &= \left\{ \mathbf{x} \in \mathbb{R}^d : \frac{|\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) - \frac{\alpha}{\beta} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}; \mathbf{x})|}{\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})} \leq T \right\} \\ &= \{ \mathbf{x} \in \mathbb{R}^d : |1 - \exp(p(\mathbf{x}))| \leq T \}, \end{aligned}$$

where $p(\mathbf{x}) = -\frac{1}{2}(\boldsymbol{\mu} - \mathbf{x})^T \boldsymbol{\Lambda}^{-1}(\boldsymbol{\mu} - \mathbf{x}) + \frac{1}{2} \mathbf{x}^T \mathbf{x} + \log(\alpha/(\sqrt{|\boldsymbol{\Lambda}|}\beta))$. Using Theorem 1 for the degree 2 polynomial $p(\mathbf{x})$ and setting $q = 4$, $\gamma = \alpha^2(\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} p^2(\mathbf{x}))^{1/2}/(256C^2)$, where C is the absolute constant of Theorem 1, we get that

$$\mathcal{N}_0(\{z : |p(z)| \leq \gamma\}) \leq \frac{\alpha}{4}.$$

To simplify notation set $Q = \{z : |p(z)| \leq \gamma\}$. Therefore, for any \mathbf{x} in the remaining $\alpha/4$ mass of the set S we know that $|p(\mathbf{x})| \geq \gamma$. Next, we lower bound γ in terms of the distance of the parameters of the two Gaussians. We have

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [p^2(\mathbf{x})] &\geq \text{Var}_{\mathbf{x} \sim \mathcal{N}_0} [p(\mathbf{x})] \\ &= \text{Var}_{\mathbf{x} \sim \mathcal{N}_0} \left[-\frac{1}{2}(\boldsymbol{\mu} - \mathbf{x})^T \boldsymbol{\Lambda}^{-1}(\boldsymbol{\mu} - \mathbf{x}) + \frac{1}{2} \mathbf{x}^T \mathbf{x} \right] \\ &= \text{Var}_{\mathbf{x} \sim \mathcal{N}_0} \left[\sum_{i=1}^d \left(\frac{\mu_i}{\lambda_i} x + x^2 \frac{(1 - 1/\lambda_i)}{2} \right) \right] \\ &= \sum_{i=1}^d \text{Var}_{\mathbf{x} \sim \mathcal{N}(0,1)} \left[\frac{\mu_i}{\lambda_i} x + x^2 \frac{(1 - 1/\lambda_i)}{2} \right] \\ &= \sum_{i=1}^d \frac{1}{2} \left(\frac{1}{\lambda_i} - 1 \right)^2 + \frac{\mu_i^2}{\lambda_i^2} \\ &= \frac{1}{2} \|\boldsymbol{\Lambda}^{-1} - \mathbf{I}\|_F^2 + \|\boldsymbol{\Lambda}^{-1/2} \boldsymbol{\mu}\|_2^2 \end{aligned}$$

Therefore, using the inequality $\sqrt{2}\sqrt{x+y} \geq \sqrt{x} + \sqrt{y}$ we obtain

$$\begin{aligned} \gamma &\geq \frac{\alpha^2}{256\sqrt{2}C^2} \left(\frac{1}{\sqrt{2}} \|\boldsymbol{\Lambda}^{-1} - \mathbf{I}\|_F + \|\boldsymbol{\Lambda}^{-1/2} \boldsymbol{\mu}\|_2 \right) \\ &\geq \frac{\alpha^2}{256C^2} d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)), \end{aligned}$$

where we used Lemma 1. Assume first that $\gamma \leq 1$. We have that the L_1 distance between the functions $f(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) \mathbf{1}_S(\mathbf{x})$ and $g(\mathbf{x}) = \frac{\alpha}{\beta} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}; \mathbf{x}) \mathbf{1}_S(\mathbf{x})$ is

$$\begin{aligned} &\int |f(\mathbf{x}) - g(\mathbf{x})| d\mathbf{x} \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [|1 - \exp(p(\mathbf{x}))| \mathbf{1}_S(\mathbf{x})] \\ &\geq \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} \left[\frac{|p(\mathbf{x})|}{2} \mathbf{1}_{S \setminus Q}(\mathbf{x}) \right] \\ &\geq \gamma \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\mathbf{1}_{S \setminus Q}(\mathbf{x})] \geq \frac{\alpha\gamma}{4} \\ &\geq C_\alpha d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)), \end{aligned}$$

where for the first inequality we used the inequality $|1 - e^x| \geq |x|/2$ for $|x| \leq 1$. Note that $C_\alpha = \Omega(\alpha^3)$. If $\gamma > 1$ we have

$$\begin{aligned} \int |f(\mathbf{x}) - g(\mathbf{x})| d\mathbf{x} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [|1 - \exp(p(\mathbf{x}))| \mathbf{1}_S(\mathbf{x})] \\ &\geq \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} \left[\frac{1}{2} \mathbf{1}_{S \setminus Q}(\mathbf{x}) \right] \geq \alpha/8, \end{aligned}$$

where we used the inequality $|1 - e^x| \geq 1/2$ for $|x| > 1$. ■

A. Learning a Weighted Characteristic Function

Our goal in this section is to recover using conditional samples from \mathcal{N}_S^* a weighted characteristic function of the set S . In particular, we will show that it is possible to learn a good approximation to the function

$$\psi(\mathbf{x}) = \frac{\mathbf{1}_S(\mathbf{x}) \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})}{\alpha^* \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})} = \frac{\mathbf{1}_S(\mathbf{x}) \mathcal{N}^*(\mathbf{x})}{\alpha^* \mathcal{N}_0(\mathbf{x})}. \quad (9)$$

We will later use the knowledge of this function to extract the unknown parameters and learn the set S .

1) *Hermite Concentration:* We start by showing that the function $\psi(\mathbf{x})$ admits strong Hermite concentration. This means that we can well-approximate $\psi(\mathbf{x})$ if we ignore the higher order terms in the Hermite expansion of $\psi(\mathbf{x})$.

Theorem 3. (LOW DEGREE APPROXIMATION) *Let $S_k \psi$ denote the degree k Hermite expansion of function ψ defined in (9). We have that*

$$\begin{aligned} &\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(S_k \psi(\mathbf{x}) - \psi(\mathbf{x}))^2] \\ &= \sum_{|V| \geq k} \hat{\psi}(V)^2 \leq \text{poly}(1/\alpha) \left(\frac{\sqrt{\Gamma(S)}}{k^{1/4}} + \frac{1}{k} \right). \end{aligned}$$

where $\Gamma(S)$ is the Gaussian surface area of S , and $a < \alpha^*$ is the absolute constant of (3).

We note that the Hermite expansion of ψ is well-defined as $\psi(\mathbf{x}) \in L_2(\mathbb{R}^d, \mathcal{N}_0)$. This can be seen from the following lemma which will be useful in many calculations throughout the paper.

Lemma 9. *Let $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ be two $(B, \frac{1-\delta}{2k})$ -isotropic Gaussians for some parameters $B, \delta > 0$ and $k \in \mathbb{N}$. Let $C = \frac{13k^2}{\delta}B$. It holds*

$$\exp(-C) \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} \left[\left(\frac{\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1; \mathbf{x})}{\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2; \mathbf{x})} \right)^k \right] \leq \exp(C).$$

Lemma 9 applied for \mathcal{N}_0 and \mathcal{N}^* for $k = 2$ implies that $\psi(\mathbf{x}) \in L_2(\mathbb{R}^d, \mathcal{N}_0)$.

To get the desired bound for Theorem 3 we use the following lemma, which allows us to bound the Hermite concentration of a function f through its noise stability.

Lemma 10. *For any function $f : \mathbb{R}^d \mapsto \mathbb{R}$ and parameter $\rho \in (0, 1)$, it holds*

$$\sum_{|V| \geq 1/\rho} \hat{f}(V)^2 \leq 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [f(\mathbf{x})^2 - f(\mathbf{x})T_{1-\rho}f(\mathbf{x})]$$

Lemma 10 was originally shown in [KKMS05] for indicator functions of sets, but their proof extends to arbitrary real functions. For a proof see the full of the paper.

Using Lemma 10, we can obtain Theorem 3 by bounding the noise sensitivity of the function ψ . The following lemma directly gives the desired result.

Lemma 11. *For any $\rho \in (0, 1)$,*

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\psi(\mathbf{x})^2 - \psi(\mathbf{x})T_{1-\rho}\psi(\mathbf{x})] \\ & \leq \text{poly}(1/\alpha) \left(\sqrt{\Gamma(S)}\rho^{1/4} + \rho \right) \end{aligned}$$

To prove Lemma 11, we will require the following lemma whose proof is provided in the full version of the paper.

Lemma 12. *Let $r(\mathbf{x}) \in L_2(\mathbb{R}^d, \mathcal{N}(\mathbf{0}, \mathbf{I}))$ be differentiable at every $\mathbf{x} \in \mathbb{R}^d$. Then*

$$\frac{1}{2} \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [(r(\mathbf{x}) - r(\mathbf{z}))^2] \leq \rho \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\nabla r(\mathbf{x})\|_2^2]$$

We now move on to the proof of Lemma 11.

Proof of Lemma 11: For ease of notation we define the following distribution

$$D_\rho = \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} \mathbf{I} & (1-\rho)\mathbf{I} \\ (1-\rho)\mathbf{I} & \mathbf{I} \end{pmatrix} \right).$$

We also denote by $r(\mathbf{x}) = \mathcal{N}^*(\mathbf{x})/\mathcal{N}_0(\mathbf{x})$ We can now write

$$\begin{aligned} & 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\psi(\mathbf{x})^2 - \psi(\mathbf{x})T_{1-\rho}\psi(\mathbf{x})] \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [\psi(\mathbf{x})^2 - \psi(\mathbf{x})\psi(\mathbf{z})] \\ & = \frac{1}{\alpha^{*2}} \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [\mathbf{1}_S(\mathbf{x})r^2(\mathbf{x}) - \mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z})r^2(\mathbf{x})] + \\ & \quad \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [\mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z})r^2(\mathbf{x}) - \mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z})r(\mathbf{x})r(\mathbf{z})] \end{aligned}$$

We bound each of the two terms separately. For the first term, using Schwarz's inequality we get

$$\begin{aligned} & \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [\mathbf{1}_S(\mathbf{x})r^2(\mathbf{x}) - \mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z})r^2(\mathbf{x})] \\ & \leq \left(\mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [\mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z})] \right)^{1/2} \left(\mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [r^4(\mathbf{x})] \right)^{1/2} \\ & \leq (\text{NS}[S])^{1/2} \text{poly}(1/\alpha) \leq \sqrt{\Gamma(S)}\rho^{1/4} \text{poly}(1/\alpha) \end{aligned}$$

where the bound on the expectation of $r^4(\mathbf{x})$ follows from Lemma 9 and the last inequality follows from Lemma 4.

For the second term, we have that

$$\begin{aligned} & \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [\mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z})(r^2(\mathbf{x}) - r(\mathbf{x})r(\mathbf{z}))] \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} \left[\mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z}) \left(\frac{r^2(\mathbf{x})}{2} + \frac{r^2(\mathbf{z})}{2} - r(\mathbf{x})r(\mathbf{z}) \right) \right] \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} \left[\mathbf{1}_S(\mathbf{x})\mathbf{1}_S(\mathbf{z}) \frac{1}{2} (r(\mathbf{x}) - r(\mathbf{z}))^2 \right] \\ & \leq \frac{1}{2} \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim D_\rho} [(r(\mathbf{x}) - r(\mathbf{z}))^2] \leq \rho \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\|\nabla r(\mathbf{x})\|_2^2], \end{aligned}$$

where the last inequality follows from Lemma 12. It thus suffices to bound the expectation of the gradient of r . We have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\|\nabla r(\mathbf{x})\|_2^2] \\ & = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} \left[\left\| -\boldsymbol{\Sigma}^{*-1}(\mathbf{x} - \boldsymbol{\mu}^*) + \mathbf{x} \right\|_2^2 r^2(\mathbf{x}) \right] \\ & \leq 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\|(\mathbf{I} - \boldsymbol{\Sigma}^*)^{-1}\mathbf{x}\|_2^2 r^2(\mathbf{x})] \\ & \quad + 2 \left\| \boldsymbol{\Sigma}^{*-1}\boldsymbol{\mu}^* \right\|_2^2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [r^2(\mathbf{x})] \\ & \leq 2 \sqrt{\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\|(\mathbf{I} - \boldsymbol{\Sigma}^{*-1})\mathbf{x}\|_2^4]} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [r^4(\mathbf{x})] \\ & \quad + 2 \left\| \boldsymbol{\Sigma}^{*-1}\boldsymbol{\mu}^* \right\|_2^2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [r^2(\mathbf{x})] \leq \text{poly}(1/\alpha) \end{aligned}$$

where the bound on the expectation of $r^4(\mathbf{x})$ and $r^2(\mathbf{x})$ follows from Lemma 9 and the expectation

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [\|(\mathbf{I} - \boldsymbol{\Sigma}^{*-1})\mathbf{x}\|_2^4] = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} \left[\left(\sum_i (1 - \lambda_i)^2 x_i^2 \right)^2 \right] \\ & \leq 3 \left(\sum_i (1 - \lambda_i)^2 \right)^2 \leq 3 \log^2(1/\alpha) \leq \text{poly}(1/\alpha) \end{aligned}$$

□

2) *Learning the Hermite Expansion:* In this section we deal with the sample complexity of estimating the coefficients of the Hermite expansion. We have

$$c_V = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)} [H_V(\mathbf{x})]$$

Using samples \mathbf{x}_i from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$, we can estimate this expectation empirically with the unbiased estimate

$$\tilde{c}_V = \frac{\sum_{i=1}^N H_V(\mathbf{x}_i)}{N}.$$

We now show an upper bound for the variance of the above estimate. The proof of this lemma can be found in the full version of the paper.

Lemma 13. *Let $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$ be the unknown truncated Gaussian. The variance of the following unbiased estimator of the Hermite coefficients $\tilde{c}_V = \frac{\sum_{i=1}^N H_V(\mathbf{x}_i)}{N}$, is upper bounded*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)} [(\tilde{c}_V - c_V)^2] \leq \text{poly}(1/\alpha) \frac{5^{|V|}}{N}.$$

Theorem 4. *Let S be an arbitrary (Borel) subset of \mathbb{R}^d . Let α be the constant of (3). Let $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$ be the corresponding truncated Gaussian in $(O \log(1/\alpha), 1/16)$ -isotropic position (see Definition 6), Then, for the estimate*

$$\psi_k(\mathbf{x}) = \max \left(0, \sum_{V: 0 \leq |V| \leq k} \tilde{c}_V H_V(\mathbf{x}) \right),$$

where $\tilde{c}_V = \frac{\sum_{i=1}^N H_V(\mathbf{x}_i)}{N}$ it holds for $k \ll d$, $\Gamma(S) > 1$,

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)} \left[\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\psi_k(\mathbf{x}) - \psi(\mathbf{x}))^2] \right] \\ & \leq \text{poly}(1/\alpha) \left(\frac{\sqrt{\Gamma(S)}}{k^{1/4}} + \frac{(5d)^k}{N} \right). \end{aligned}$$

Alternatively, for $k = \text{poly}(1/\alpha) \Gamma(S)^2 / \varepsilon^4$ we obtain that with $N = d^{\text{poly}(1/\alpha) \Gamma(S)^2 / \varepsilon^4}$ samples, with probability at least $9/10$, it holds $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(p_{N,k}(\mathbf{x}) - \psi(\mathbf{x}))^2] \leq \varepsilon$.

Proof: Instead of considering the positive part of the Hermite expansion, we will prove the claim for the empirical Hermite expansion of degree k and N samples

$$p_{N,k} = \sum_{V: 0 \leq |V| \leq k} \tilde{c}_V H_V(\mathbf{x}).$$

As usual we denote by $S_k \psi(\mathbf{x})$ the true (exact) Hermite expansion of degree k of $\psi(\mathbf{x})$. Using the inequality $(a - b)^2 \leq 2(a - c)^2 + 2(c - b)^2$ we obtain

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(p_{N,k}(\mathbf{x}) - f(\mathbf{x}))^2] \\ & \leq 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(p_{N,k}(\mathbf{x}) - S_k \psi(\mathbf{x}))^2] \\ & \quad + 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(S_k \psi(\mathbf{x}) - \psi(\mathbf{x}))^2] \end{aligned}$$

Since Hermite polynomials form an orthonormal system with respect to \mathcal{N}_0 , we obtain

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(p_{N,k}(\mathbf{x}) - S_k \psi(\mathbf{x}))^2] = \sum_{V: 0 \leq |V| \leq k} (\tilde{c}_V - c_V)^2.$$

Using Lemma 13 we obtain

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathcal{N}^*} \left[\sum_{V: 0 \leq |V| \leq k} (\tilde{c}_V - c_V)^2 \right] \\ & \leq \frac{\text{poly}(1/\alpha)}{N} \sum_{V: 0 \leq |V| \leq k} 5^{|V|} \leq \frac{\text{poly}(1/\alpha)}{N} \binom{d+k}{k} 5^k, \end{aligned}$$

where we used the fact that the number of all multi-indices V of d elements such that $0 \leq |V| \leq k$ is $\binom{d+k}{k}$. Moreover, from Theorem 3 we obtain that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(S_k \psi(\mathbf{x}) - \psi(\mathbf{x}))^2] \leq \text{poly}(1/\alpha) \left(\frac{\sqrt{\Gamma(S)}}{k^{1/4}} + \frac{1}{k} \right).$$

The theorem follows. \blacksquare

IV. ESTIMATION ALGORITHM FOR BOUNDED GAUSSIAN SURFACE AREA

In this section, we present the main steps of our estimation algorithm. In later sections, we provide details of the individual components. The algorithm can be thought of in 3 stages.

First Stage: In the first stage, our goal is to learn a weighted characteristic function of the underlying set. Even though we cannot access the underlying set directly, for any given function f we can evaluate the expectation $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)} [f(\mathbf{x})]$ using truncated samples.

This expectation can be equivalently written as $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [f(\mathbf{x}) \psi(\mathbf{x})]$ for the function

$$\psi(\mathbf{x}) := \frac{\mathbf{1}_S(\mathbf{x}) \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})}{\alpha^* \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})} = \frac{\mathbf{1}_S(\mathbf{x}) \mathcal{N}^*(\mathbf{x})}{\alpha^* \mathcal{N}_0(\mathbf{x})}.$$

By evaluating the above expectation for different functions f corresponding to the Hermite polynomials $H_V(\mathbf{x})$, we can recover $\psi(\mathbf{x})$, through its Hermite expansion:

$$\begin{aligned} \psi(\mathbf{x}) &= \sum_{V \in \mathbb{N}^d} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [H_V(\mathbf{x}) \psi(\mathbf{x})] H_V(\mathbf{x}) \\ &= \sum_{V \in \mathbb{N}^d} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} [H_V(\mathbf{x})] H_V(\mathbf{x}). \end{aligned}$$

Of course, it is infeasible to calculate the Hermite expansion for any $V \in \mathbb{N}^d$. In Section III-A, we show that by estimating only terms of degree at most k , we can achieve a good approximation to ψ where the error depends on the Gaussian surface area of the underlying set S . To do this, we show that most of the mass of the coefficients $c_V = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [H_V(\mathbf{x}) \psi(\mathbf{x})]$ is concentrated on low degree terms, i.e. $\sum_{|V| > k} c_V^2$ is significantly small. Moreover, we show that even though we can only estimate the coefficients c_V through sampling, the sampling error is significantly small.

Overall, after the first stage, we obtain a non-negative function ψ_k that is close to ψ . The approximation error guarantees are given in Theorem 4.

Second Stage: Given the function ψ_k that was recovered in the first stage, our goal is to decouple the influence of the set $\frac{1_S(\mathbf{x})}{\alpha^*}$ and the influence of the underlying Gaussian distribution which corresponds to the multiplicative term $\frac{\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})}{\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})}$. This would be easy if we had the exact function ψ in hand. In contrast, for the polynomial function ψ_k the problem is significantly challenging as it is only close to ψ on average but not pointwise.

To perform the decoupling and identify the underlying Gaussian we explicitly multiply the function ψ_k with a corrective term of the form $\frac{\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})}{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{x})}$. We set up an optimization problem seeking to minimize the function $C(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*}[\frac{\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})}{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{x})} \psi_k(\mathbf{x})]$ with an appropriate choice of $C(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ so that the unique solution corresponds to $(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$. Under a reparameterization of $(\mathbf{u}, \mathbf{B}) = (\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1})$, we show that the corresponding problem is *strongly* convex. Still, optimizing it directly is non-trivial as it involves taking the expectation with respect to the unknown truncated Gaussian. Instead, we perform stochastic gradient descent (SGD) and show that it quickly converges in few steps to point close to the true minimizer (Algorithm 2).

This allows us to recover parameters $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$ so that the total variation distance between the recovered and the true (untruncated) Gaussian is very small, i.e. $d_{TV}(\mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}), \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)) \leq \varepsilon$. Theorem 5 describes the guarantees of the second stage. Further details are provided in Section IV-A.

The guarantees of the algorithm: We first show our algorithmic results under the assumption that the untruncated Gaussian \mathcal{N}^* is known to be in near-isotropic position.

Definition 6 (Near-Isotropic Position). *Let $\boldsymbol{\mu} \in \mathbb{R}^d$, $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ be a positive semidefinite symmetric matrix and $a, b > 0$. We say that $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is in (a, b) -isotropic position if the following hold.*

$$\|\boldsymbol{\mu}\|_2^2 \leq a, \quad \|\boldsymbol{\Sigma} - \mathbf{I}\|_F^2 \leq a, \quad (1-b)\mathbf{I} \preceq \boldsymbol{\Sigma} \preceq \frac{1}{1-b}\mathbf{I}$$

We later transform the more interesting case with an unknown mean and an unknown diagonal covariance matrix to the isotropic case.

Theorem 5. *Let $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$ be a d -dimensional Gaussian distribution that is in $(O(\log(1/\alpha^*)), 1/16)$ -isotropic position and consider a set S such that $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; S) \geq \alpha$. There exists an algorithm such that for all $\varepsilon > 0$, the algorithm uses $n > d^{\text{poly}(1/\alpha)} \frac{\Gamma^2(S)}{\varepsilon^8}$ samples and produces, in $\text{poly}(n)$ time, estimates that, with probability at least 99%, satisfy $d_{TV}(\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*), \mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})) \leq \varepsilon$.*

We can apply this theorem to estimate the parameters of any Gaussian distribution with an unknown mean and an unknown diagonal covariance matrix by bringing the Gaussian to an $(O(\log(1/\alpha^*)), 1/16)$ -isotropic position. Lemma 3

shows that with high probability, we can obtain initial estimates $\tilde{\boldsymbol{\mu}}_S$ and $\tilde{\boldsymbol{\Sigma}}_S$ so that $\|\boldsymbol{\Sigma}^{-1/2}(\tilde{\boldsymbol{\mu}}_S - \boldsymbol{\mu}^*)\|_2^2 \leq O(\log \frac{1}{\alpha})$ and

$$\tilde{\boldsymbol{\Sigma}}_S \succeq \Omega(\alpha^2) \boldsymbol{\Sigma}^*,$$

$$\text{and } \left\| \boldsymbol{\Sigma}^{*-1/2} \tilde{\boldsymbol{\Sigma}}_S \boldsymbol{\Sigma}^{*-1/2} - \mathbf{I} \right\|_F^2 \leq O(\log \frac{1}{\alpha}).$$

Given these estimates, we can transform the space so that $\tilde{\boldsymbol{\mu}}_S = \mathbf{0}$, and $\tilde{\boldsymbol{\Sigma}}_S = \mathbf{I}$. We note that after this transformation, the mean will be at the right distance from 0, while the eigenvalues λ_i of $\boldsymbol{\Sigma}^*$ will all be within the desired range $\frac{15}{16} \leq \lambda_i \leq \frac{16}{15}$ apart from at most $O(\log(1/\alpha))$. This is because the condition $\left\| \boldsymbol{\Sigma}^{*-1/2} \tilde{\boldsymbol{\Sigma}}_S \boldsymbol{\Sigma}^{*-1/2} - \mathbf{I} \right\|_F^2 \leq O(\log \frac{1}{\alpha})$ implies that $\sum_i (1 - \frac{1}{\lambda_i})^2 \leq O(\log(1/\alpha))$. With this observation, since we know of the eigenvectors of $\boldsymbol{\Sigma}^*$, we would be able to search over all possible corrections to the eigenvalues to bring the Gaussian in $(O(\log(1/\alpha)), \frac{1}{16})$ -isotropic position as required by Theorem 5. We only need to correct $O(\log(1/\alpha))$ of them.

We can form a space of candidate hypotheses for the underlying distribution, for each choice of $O(\log(1/\alpha))$ out of the d vectors along with the all possible scalings. These hypotheses are at most $d^{O(\log(1/\alpha))}$ times $(\log(1/\alpha))^{O(\log(1/\alpha))}$ for all possible scalings. Thus, there are at most $d^{O(\log(1/\alpha))}$ hypotheses. Running the algorithm for each one of them, we would learn at least one distribution and one set that is accurate according to the guarantees of Theorems 5. Running the generic hypothesis testing algorithm of Lemma 6, we can identify one that is closest in total variation distance to the true distribution \mathcal{N}_S^* . The sample complexity and runtime would thus only increase by at most $d^{O(\log(1/\alpha))}$. As we showed in Lemma 8, knowing the truncated Gaussian in total variation distance suffices to learn in accuracy ε the parameters of the untruncated distribution. We thus obtain as corollary, that we can estimate the parameters when the covariance is spherical or diagonal. The same results hold when one wants to recover the underlying set in these cases.

A. Optimization of Gaussian Parameters

In this section we show that we can formulate a convex objective function that can be optimized to yield the unknown parameters $\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*$ of the truncated Gaussian. Let S be the unknown (Borel) subset of \mathbb{R}^d such that $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; S) = \alpha^*$ and let $\mathcal{N}_S^* = \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$ be the corresponding truncated Gaussian.

To find the parameters $\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*$, we define the function

$$M_f(\mathbf{u}, \mathbf{B}) := \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) f(\mathbf{x}) \right] \quad (10)$$

where $h(\mathbf{u}, \mathbf{B}; \mathbf{x}) = \frac{\mathbf{x}^T \mathbf{B} \mathbf{x}}{2} - \frac{\text{tr}((\mathbf{B} - \mathbf{I})(\tilde{\boldsymbol{\Sigma}}_S + \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T))}{2} - \mathbf{u}^T (\mathbf{x} - \tilde{\boldsymbol{\mu}}_S) + \frac{d}{2} \log 2\pi$.

We will show that the minimizer of $M_f(\mathbf{u}, \mathbf{B})$ for the polynomial function $f = \psi_k$, will satisfy $(\mathbf{B}^{-1} \mathbf{u}, \mathbf{B}^{-1}) \approx$

$(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$. Note that $M_f(\mathbf{u}, \mathbf{B})$ can be estimated through samples. Our goal will be to optimize it through stochastic gradient descent.

In order to make sure that SGD algorithm for M_{ψ_k} converges fast in the parameter space we need to project after every iteration to some subset of the space as we will see in more details later in this Section. Assuming that the pair $(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$ is in $(\sqrt{\log(1/\alpha^*)}, 1/16)$ -isotropic position we define the following set

$$\mathcal{D} = \{(\mathbf{u}, \mathbf{B}) \mid (\mathbf{B}^{-1}\mathbf{u}, \mathbf{B}^{-1}) \text{ is in } (c \cdot \log(1/\alpha^*), 1/16)\text{-isotropic position}\} \quad (11)$$

Where c is the universal constant guaranteed to exist from Section II so that

$$\max \left\{ \|\boldsymbol{\mu}^* - \tilde{\boldsymbol{\mu}}\|_{\boldsymbol{\Sigma}^*}, \|\boldsymbol{\Sigma}^* - \tilde{\boldsymbol{\Sigma}}\|_F \right\} \leq c \cdot \log(1/\alpha^*).$$

It is not hard to see that \mathcal{D} is a convex set and that for any (\mathbf{u}, \mathbf{B}) the projection to \mathcal{D} can be done efficiently. For more details we refer to Lemma 8 of [DGTZ18]. Since after every iteration of our algorithm we project to \mathcal{D} we will assume for the rest of this Section that $(\mathbf{u}, \mathbf{B}) \in \mathcal{D}$.

An equivalent formulation of $M_f(\mathbf{u}, \mathbf{B})$ that will be useful for the analysis of the SGD algorithm is

$$\begin{aligned} M_f(\mathbf{u}, \mathbf{B}) &= e^{-\frac{1}{2}(\text{tr}((\mathbf{B}-\mathbf{I})(\tilde{\boldsymbol{\Sigma}}_S + \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T)) + \mathbf{u}^T \mathbf{B}^{-1} \mathbf{u} - \mathbf{u}^T \tilde{\boldsymbol{\mu}}_S)} \sqrt{|\mathbf{B}|} \\ &\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[\frac{\mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})}{\mathcal{N}(\mathbf{B}^{-1}\mathbf{u}, \mathbf{B}^{-1}; \mathbf{x})} f(\mathbf{x}) \right] \\ &:= C_{\mathbf{u}, \mathbf{B}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[\frac{\mathcal{N}_0(\mathbf{x})}{\mathcal{N}_{\mathbf{u}, \mathbf{B}}(\mathbf{x})} f(\mathbf{x}) \right] \end{aligned} \quad (12)$$

Lemma 14. For $(\mathbf{u}, \mathbf{B}) \in \mathcal{D}$, we have that $\text{poly}(\alpha) \leq C_{\mathbf{u}, \mathbf{B}} \leq \text{poly}(1/\alpha)$.

Proof: We have that

$$\begin{aligned} |2 \log C_{\mathbf{u}, \mathbf{B}}| &= |\text{tr}((\mathbf{B}-\mathbf{I})(\tilde{\boldsymbol{\Sigma}}_S + \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T)) \\ &\quad + \mathbf{u}^T \mathbf{B}^{-1} \mathbf{u} - \mathbf{u}^T \tilde{\boldsymbol{\mu}}_S - \log |\mathbf{B}|)| \\ &= |\text{tr}(\mathbf{B}-\mathbf{I}) + \text{tr}((\mathbf{B}-\mathbf{I})(\tilde{\boldsymbol{\Sigma}}_S - \mathbf{I})) \\ &\quad + \mathbf{u}^T \mathbf{B}^{-1} \mathbf{u} - \log |\mathbf{B}| \\ &\leq |\text{tr}(\mathbf{B}-\mathbf{I}) - \log |\mathbf{B}|| \\ &\quad + |\text{tr}((\mathbf{B}-\mathbf{I})(\tilde{\boldsymbol{\Sigma}}_S - \mathbf{I}))| + |\mathbf{u}^T \mathbf{B}^{-1} \mathbf{u}| \end{aligned}$$

We now bound each of the terms separately. Let $\lambda_1, \dots, \lambda_d$ be the eigenvalues of \mathbf{B} .

For the first term, we have that

$$\begin{aligned} |\text{tr}(\mathbf{B}-\mathbf{I}) - \log |\mathbf{B}|| &= \left| \sum_{i=1}^d (\lambda_i - 1 - \log \lambda_i) \right| \leq \sum_{i=1}^d \frac{(\lambda_i - 1)^2}{\lambda_i} \leq \frac{\|\mathbf{B}-\mathbf{I}\|_F^2}{\lambda_{\min}}, \end{aligned}$$

where we used the fact that $0 \leq x - 1 - \log x \leq \frac{(x-1)^2}{x}$ for all $x > 0$. For the second term, we have that $|\text{tr}((\mathbf{B}-\mathbf{I})(\tilde{\boldsymbol{\Sigma}}_S - \mathbf{I}))| \leq \|\mathbf{B}-\mathbf{I}\|_F \|\tilde{\boldsymbol{\Sigma}}_S - \mathbf{I}\|_F$. For the third term, we have that $|\mathbf{u}^T \mathbf{B}^{-1} \mathbf{u}| = \mathbf{u}^T \mathbf{B}^{-1} \mathbf{B} \mathbf{B}^{-1} \mathbf{u} \leq \lambda_{\max} \|\mathbf{B}^{-1} \mathbf{u}\|_2^2$.

Now from the assumption $(\mathbf{u}, \mathbf{B}) \in \mathcal{D}$ we have that $\|\mathbf{B}-\mathbf{I}\|_F \leq O(\sqrt{\log(1/\alpha^*)})$, $\|\mathbf{B}^{-1} \mathbf{u}\|_2 \leq O(\sqrt{\log(1/\alpha^*)})$, $\lambda_{\min} \geq 15/16$ and $\lambda_{\max} \leq 17/16$. Also from Lemma 3 we get that $\|\tilde{\boldsymbol{\Sigma}}_S - \mathbf{I}\|_F \leq O(\sqrt{\log(1/\alpha^*)})$ and hence $|2 \log C_{\mathbf{u}, \mathbf{B}}| \leq O(\log(1/\alpha^*))$. This means that $C_{\mathbf{u}, \mathbf{B}} = \text{poly}(1/\alpha)$ and the lemma follows. \blacksquare

1) *The Objective Function and its Approximation:*

To show that the minimizer of the function M_{ψ_k} is a good estimator for the unknown parameters $\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*$, we consider the function M'_f , defined as $M'_f(\mathbf{u}, \mathbf{B}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[e^{h'(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) f(\mathbf{x}) \right]$ for $h'(\mathbf{u}, \mathbf{B}; \mathbf{x}) = \frac{\mathbf{x}^T \mathbf{B} \mathbf{x}}{2} - \frac{\text{tr}((\mathbf{B}-\mathbf{I})(\boldsymbol{\Sigma}_S + \boldsymbol{\mu}_S \boldsymbol{\mu}_S^T))}{2} - \mathbf{u}^T (\mathbf{x} - \boldsymbol{\mu}_S) + \frac{d}{2} \log 2\pi$. This function corresponds to an ideal situation where we know the parameters $\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S$ exactly. Similarly to (12), we can write M'_f as $C'_{\mathbf{u}, \mathbf{B}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[\frac{\mathcal{N}_0(\mathbf{x})}{\mathcal{N}_{\mathbf{u}, \mathbf{B}}(\mathbf{x})} f(\mathbf{x}) \right]$. We argue that both M_f and M'_f are convex.

Claim 1. For any function $f : \mathbb{R}^d \mapsto \mathbb{R}_{\geq 0}$, $M_f(\mathbf{u}, \mathbf{B})$ and $M'_f(\mathbf{u}, \mathbf{B})$ are convex functions of the parameters (\mathbf{u}, \mathbf{B}) .

Proof: We show the statement for M_f . The proof for M'_f is identical. The proof follows by computing the Hessian of M_f and arguing that it is positive semidefinite.

The gradient with respect to (\mathbf{u}, \mathbf{B}) is

$$\begin{aligned} \nabla M_f(\mathbf{u}, \mathbf{B}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)} \left[\nabla h(\mathbf{u}, \mathbf{B}; \mathbf{x}) e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) f(\mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)} \left[\left(\frac{1}{2} \left(\mathbf{x} \mathbf{x}^T - \tilde{\boldsymbol{\Sigma}}_S - \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T \right)^b \right) \right. \\ &\quad \left. \frac{\tilde{\boldsymbol{\mu}}_S - \mathbf{x}}{e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) f(\mathbf{x})} \right] \end{aligned} \quad (13)$$

Moreover, the Hessian is

$$\begin{aligned} \mathcal{H}_{M_f}(\mathbf{u}, \mathbf{B}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)} \left[\left(\frac{1}{2} \left(\mathbf{x} \mathbf{x}^T - \tilde{\boldsymbol{\Sigma}}_S - \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T \right)^b \right) \right. \\ &\quad \left. \left(\frac{1}{2} \left(\mathbf{x} \mathbf{x}^T - \tilde{\boldsymbol{\Sigma}}_S - \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T \right)^b \right)^T e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) f(\mathbf{x}) \right] \end{aligned}$$

which is clearly positive semidefinite since for any $\mathbf{z} \in \mathbb{R}^{d \times d+d}$ we have

$$\begin{aligned} \mathbf{z}^T \mathcal{H}_{M_f}(\mathbf{u}, \mathbf{B}) \mathbf{z} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)} \left[\left(\mathbf{z}^T \left(\frac{1}{2} \left(\mathbf{x} \mathbf{x}^T - \tilde{\boldsymbol{\Sigma}}_S - \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T \right)^b \right) \right)^2 \right] \end{aligned}$$

$$\left. e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) f(\mathbf{x}) \right] \geq 0.$$

■

We now argue that the minimizer of the convex function M'_ψ for the weighted characteristic function $\psi(\mathbf{x}) = \frac{\mathbf{1}_S(\mathbf{x}) \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})}{\alpha^* \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})}$ is $(\mathbf{u}, \mathbf{B}) = (\boldsymbol{\Sigma}^{*-1}, \boldsymbol{\Sigma}^{*-1} \boldsymbol{\mu}^*)$.

Claim 2. *The minimizer of $M'_\psi(\mathbf{u}, \mathbf{B})$ is $(\mathbf{u}, \mathbf{B}) = (\boldsymbol{\Sigma}^{*-1}, \boldsymbol{\Sigma}^{*-1} \boldsymbol{\mu}^*)$.*

Proof: The gradient of M'_ψ with respect to (\mathbf{u}, \mathbf{B}) is

$$\begin{aligned} \nabla M'_\psi(\mathbf{u}, \mathbf{B}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[\begin{pmatrix} \frac{1}{2} (\mathbf{x}\mathbf{x}^T - \boldsymbol{\Sigma}_S - \boldsymbol{\mu}_S \boldsymbol{\mu}_S^T)^b \\ \boldsymbol{\mu}_S - \mathbf{x} \end{pmatrix} \right. \\ &\quad \left. e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) \frac{\mathbf{1}_S(\mathbf{x}) \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})}{\alpha^* \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x})} \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[\begin{pmatrix} \frac{1}{2} (\mathbf{x}\mathbf{x}^T - \boldsymbol{\Sigma}_S - \boldsymbol{\mu}_S \boldsymbol{\mu}_S^T)^b \\ \boldsymbol{\mu}_S - \mathbf{x} \end{pmatrix} \right. \\ &\quad \left. e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \frac{\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})}{\alpha^*} \right] \end{aligned}$$

For $(\mathbf{u}, \mathbf{B}) = (\boldsymbol{\Sigma}^{*-1} \boldsymbol{\mu}^*, \boldsymbol{\Sigma}^{*-1})$, this is equal to

$$\begin{aligned} &\nabla M'_\psi(\boldsymbol{\Sigma}^{*-1} \boldsymbol{\mu}^*, \boldsymbol{\Sigma}^{*-1}) \\ &= C_{\mathbf{u}, \mathbf{B}} \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[\begin{pmatrix} \frac{1}{2} (\mathbf{x}\mathbf{x}^T - \boldsymbol{\Sigma}_S - \boldsymbol{\mu}_S \boldsymbol{\mu}_S^T)^b \\ \boldsymbol{\mu}_S - \mathbf{x} \end{pmatrix} \right. \\ &\quad \left. \frac{1}{\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})} \frac{\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; \mathbf{x})}{\alpha^*} \right] \\ &= \frac{C_{\mathbf{u}, \mathbf{B}}}{\alpha^*} \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[\begin{pmatrix} \frac{1}{2} (\mathbf{x}\mathbf{x}^T - \boldsymbol{\Sigma}_S - \boldsymbol{\mu}_S \boldsymbol{\mu}_S^T)^b \\ \boldsymbol{\mu}_S - \mathbf{x} \end{pmatrix} \right] \end{aligned}$$

where $C_{\mathbf{u}, \mathbf{B}}$ that does not depend on x . This is equal to 0 by definition of $\boldsymbol{\mu}_S$ and $\boldsymbol{\Sigma}_S$. ■

We want to show that the minimizer of M_{ψ_k} is close to that of M'_ψ . To do this, we bound the difference of the two functions pointwise. The proof of the following lemma is technical and can be found in the full version of the paper.

Lemma 15 (POINTWISE APPROXIMATION OF THE OBJECTIVE FUNCTION). *Assume that we use Lemma 2 to estimate $\tilde{\boldsymbol{\mu}}_S, \tilde{\boldsymbol{\Sigma}}_S$ with $\varepsilon = \frac{1}{\text{poly}(1/\alpha^*)} \varepsilon'$ and Theorem 4 with $\varepsilon = \frac{1}{p(1/\alpha^*)} \varepsilon'^2$ then*

$$|M_{\psi_k}(\mathbf{u}, \mathbf{B}) - M'_\psi(\mathbf{u}, \mathbf{B})| \leq \varepsilon'.$$

Now that we have established that M_{ψ_k} is a good approximation of M'_ψ we will prove that we can optimize M_{ψ_k} and get a solution that is very close to the optimal solution of M'_ψ .

2) *Optimization of the Approximate Objective Function:* Our goal in this section is to prove that using sample access to $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$ we can find the minimum of the function M_{ψ_k} defined in the previous section. First of all recall that M_{ψ_k} can be written as an expectation over $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$ in the following way

$$M_{\psi_k}(\mathbf{u}, \mathbf{B}) := \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} \left[e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) \psi_k(\mathbf{x}) \right].$$

In Section III-A we prove that we can learn the function ψ_k and hence M_{ψ_k} can be written as

$$M_{\psi_k}(\mathbf{u}, \mathbf{B}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} [m_{\psi_k}(\mathbf{u}, \mathbf{B}; \mathbf{x})]$$

where $m_{\psi_k}(\mathbf{u}, \mathbf{B}; \mathbf{x}) = e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) \psi_k(\mathbf{x})$, and for any \mathbf{u}, \mathbf{B} and \mathbf{x} we can compute $m_{\psi_k}(\mathbf{u}, \mathbf{B}; \mathbf{x})$. Since M_{ψ_k} is convex we are going to use stochastic gradient descent to find its minimum. To prove the convergence of SGD and bound the number of steps that SGD needs to converge we will use the the formulation developed in Chapter 14 of [SSBD14]. To be able to use their results we have to define for any (\mathbf{u}, \mathbf{B}) a random vector $\mathbf{v}(\mathbf{u}, \mathbf{B})$ and prove the following

UNBIASED GRADIENT ESTIMATION

$$\mathbb{E}[\mathbf{v}(\mathbf{u}, \mathbf{B})] = \nabla M_{\psi_k},$$

BOUNDED STEP VARIANCE

$$\mathbb{E}[\|\mathbf{v}(\mathbf{u}, \mathbf{B})\|_2^2] \leq \rho,$$

STRONG CONVEXITY for any $\mathbf{z} \in \mathcal{D}$ it holds

$$\mathbf{z}^T \mathcal{H}_{M_f}(\mathbf{u}, \mathbf{B}) \mathbf{z} \geq \lambda.$$

We start with the definition of the random vector \mathbf{v} . Given a sample \mathbf{x} from $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$, for any (\mathbf{u}, \mathbf{B}) we define

$$\begin{aligned} \mathbf{v}(\mathbf{u}, \mathbf{B}) &= \nabla_{\mathbf{u}, \mathbf{B}} m_{\psi_k}(\mathbf{u}, \mathbf{B}; \mathbf{x}) \\ &= \begin{pmatrix} \frac{1}{2} (\mathbf{x}\mathbf{x}^T - \tilde{\boldsymbol{\Sigma}}_S - \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T)^b \\ \tilde{\boldsymbol{\mu}}_S - \mathbf{x} \end{pmatrix} \\ &\quad \cdot e^{h(\mathbf{u}, \mathbf{B}; \mathbf{x})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}) \psi_k(\mathbf{x}) \end{aligned} \quad (14)$$

observe that the randomness of \mathbf{v} only comes from the random sample $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$. The fact that $\mathbf{v}(\mathbf{u}, \mathbf{B})$ is an unbiased estimator of $\nabla M_f(\mathbf{u}, \mathbf{B})$ follows directly from the fact calculation of $\nabla M_f(\mathbf{u}, \mathbf{B})$ in Section IV-A1. For the other two properties that we need we have the following lemmas. The following lemma bounds the variance of the step of the SGD algorithm. It's rather technical proof can be found in the full version of the paper.

Lemma 16 (BOUNDED STEP VARIANCE). *Let α be the constant of (3). For every $(\mathbf{u}, \mathbf{B}) \in \mathcal{D}$ it holds*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_S^*} [\|\mathbf{v}(\mathbf{u}, \mathbf{B})\|_2^2] \leq \text{poly}(1/\alpha) \cdot d^{2k},$$

We are now going to prove the strong convexity of the objective function M_{ψ_k} . For this we are going to use the anti-concentration result, Theorem 1) for polynomial functions over the Gaussian measure.

The following lemma shows that our objective is strongly convex as long as the guess \mathbf{u}, \mathbf{B} remains in the set \mathcal{D} . For the proof, we refer the reader to the full version of the paper.

Lemma 17 (STRONG CONVEXITY). *Let α be the absolute constant of (3). For every $(\mathbf{u}, \mathbf{B}) \in \mathcal{D}$, any $\mathbf{z} \in \mathbb{R}^d$ such that $\|\mathbf{z}\|_2 = 1$ and the first d^2 coordinated of \mathbf{z} correspond to a symmetric matrix, then*

$$\mathbf{z}^T \mathcal{H}_{M_f}(\mathbf{u}, \mathbf{B}) \mathbf{z} \geq \text{poly}(\alpha),$$

3) *Recovering the Unconditional Mean and Covariance:*

The framework that we use for proving the fast convergence of our SGD algorithm is summarized in the following theorem and the following lemma.

Theorem 6 (Theorem 14.11 of [SSBD14]). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Assume that f is λ -strongly convex, that $\mathbb{E}[\mathbf{v}^{(i)} \mid \mathbf{w}^{(i-1)}] \in \partial f(\mathbf{w}^{(i-1)})$ and that $\mathbb{E}[\|\mathbf{v}^{(i)}\|_2^2] \leq \rho^2$. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathcal{D}} f(\mathbf{w})$ be an optimal solution. Then,*

$$\mathbb{E}[f(\bar{\mathbf{w}})] - f(\mathbf{w}^*) \leq \frac{\rho^2}{2\lambda T} (1 + \log T),$$

where $\bar{\mathbf{w}}$ is the output projected stochastic gradient descent with steps $\mathbf{v}^{(i)}$ and projection set \mathcal{D} after T iterations.

Lemma 18 (Lemma 13.5 of [SSBD14]). *If f is λ -strongly convex and \mathbf{w}^* is a minimizer of f , then, for any \mathbf{w} it holds that*

$$f(\mathbf{w}) - f(\mathbf{w}^*) \geq \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2.$$

Now we have all the ingredients to present the proof of Theorem 5.

Proof of Theorem 5: The estimation procedure starts by computing the polynomial function ψ_k using $d^{\text{poly}(1/\alpha^*)} \frac{r^2(S)}{\varepsilon^{1/8}}$ samples from $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$ as explained in Theorem 4 to get error $\text{poly}(\alpha^*)\varepsilon'^2$. Then we compute $\tilde{\boldsymbol{\mu}}_S$ and $\tilde{\boldsymbol{\Sigma}}_S$ as explained in Section II with $\varepsilon = \frac{q(\alpha^*)}{8p(1/\alpha^*)}(\varepsilon')^2$ where p comes from Lemma 15 and q comes from Lemma 17. Our estimators for $\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}}$ are the outputs of Algorithm 2.

We analyze the accuracy of our estimation by proving that the minimum of M_{ψ_k} is close in the parameter space to the minimum of M'_{ψ} . Let \mathbf{u}', \mathbf{B}' be the minimum of the convex function M'_{ψ} and $\mathbf{u}_k, \mathbf{B}_k$ be the minimum of the convex function M_{ψ_k} . Using Lemma 15 we have the following relations

$$\begin{aligned} |M'_{\psi}(\mathbf{u}', \mathbf{B}') - M_{\psi_k}(\mathbf{u}', \mathbf{B}')| &\leq \varepsilon', \\ |M'_{\psi}(\mathbf{u}_k, \mathbf{B}_k) - M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k)| &\leq \varepsilon' \end{aligned}$$

and also

$$\begin{aligned} M'_{\psi}(\mathbf{u}', \mathbf{B}') &\leq M'_{\psi}(\mathbf{u}_k, \mathbf{B}_k), \\ M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k) &\leq M_{\psi_k}(\mathbf{u}', \mathbf{B}'). \end{aligned}$$

These relations imply that

$$\begin{aligned} &|M_{\psi_k}(\mathbf{u}', \mathbf{B}') - M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k)| \\ &= M_{\psi_k}(\mathbf{u}', \mathbf{B}') - M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k) \\ &\leq M_{\psi_k}(\mathbf{u}', \mathbf{B}') - M'_{\psi}(\mathbf{u}', \mathbf{B}') \\ &\quad + M'_{\psi}(\mathbf{u}_k, \mathbf{B}_k) - M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k) \\ &\leq |M'_{\psi}(\mathbf{u}', \mathbf{B}') - M_{\psi_k}(\mathbf{u}', \mathbf{B}')| \\ &\quad + |M'_{\psi}(\mathbf{u}_k, \mathbf{B}_k) - M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k)| \leq 2\varepsilon'. \end{aligned}$$

But from Lemma 17 and Lemma 18 we get that $\left\| \begin{pmatrix} \mathbf{B}'^b \\ \mathbf{u}' \end{pmatrix} - \begin{pmatrix} \mathbf{B}_k^b \\ \mathbf{u}_k \end{pmatrix} \right\|_2 \leq \frac{\varepsilon'}{2}$. Now we can apply the Claim 2 which implies that

$$\left\| \begin{pmatrix} (\boldsymbol{\Sigma}^{*-1})^b \\ \boldsymbol{\Sigma}^{*-1} \boldsymbol{\mu}^* \end{pmatrix} - \begin{pmatrix} \mathbf{B}_k^b \\ \mathbf{u}_k \end{pmatrix} \right\|_2 \leq \frac{\varepsilon'}{2}. \quad (16)$$

Therefore it suffices to find $(\mathbf{u}_k, \mathbf{B}_k)$ with accuracy $\varepsilon'/2$ to get our theorem.

Let $\mathbf{w}^* = \begin{pmatrix} \mathbf{B}_k^b \\ \mathbf{u}_k \end{pmatrix}$. To prove that Algorithm 2 converges to \mathbf{w}^* we use Theorem 6 which together with Markov's inequality, Lemma 16 and Lemma 17 gives us

$$\begin{aligned} &\mathbb{P} \left(M_{\psi_k}(\hat{\mathbf{u}}, \hat{\mathbf{B}}) - M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k) \right. \\ &\quad \left. \geq \text{poly}(1/\alpha^*) \cdot \frac{d^{2k}}{T} (1 + \log(T)) \right) \leq \frac{1}{3}. \quad (17) \end{aligned}$$

To get our estimation we first repeat the SGD procedure $K = \log(1/\delta)$ times independently, with parameters T, λ each time. We then get the set of estimates $\mathcal{E} = \{\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2, \dots, \bar{\mathbf{w}}_K\}$. Because of (17) we know that, with high probability $1 - \delta$, for at least the $2/3$ of the points $\bar{\mathbf{w}}$ in \mathcal{E} it is true that $M_{\psi_k}(\mathbf{w}) - M_{\psi_k}(\mathbf{w}^*) \leq \eta$ where $\eta = \text{poly}(1/\alpha^*) \cdot \frac{d^{2k}}{T} (1 + \log(T))$. Moreover we will prove later that $M_{\psi_k}(\mathbf{w}) - M_{\psi_k}(\mathbf{w}^*) \leq \eta$ and this implies $\|\mathbf{w} - \mathbf{w}^*\| \leq c \cdot \eta$, where c is a universal constant. Therefore with high probability $1 - \delta$ for at least the $2/3$ of the points $\bar{\mathbf{w}}, \bar{\mathbf{w}}'$ in \mathcal{E} it is true that $\|\bar{\mathbf{w}} - \bar{\mathbf{w}}'\| \leq 2c \cdot \eta$. Hence if we set $\hat{\mathbf{w}}$ to be a point that is at least $2c \cdot \eta$ close to more that the half of the points in \mathcal{E} then with high probability $1 - \delta$ we have that $f(\bar{\mathbf{w}}) - f(\mathbf{w}^*) \leq \eta$. Hence we can we lose probability at most δ if we condition on the event

$$\begin{aligned} &M_{\psi_k}(\hat{\mathbf{u}}, \hat{\mathbf{B}}) - M_{\psi_k}(\mathbf{u}_k, \mathbf{B}_k) \\ &\leq \text{poly}(1/\alpha^*) \cdot \frac{d^{2k}}{T} (1 + \log(T)). \end{aligned}$$

Using once again Lemma 18 we get that

$$\left\| \begin{pmatrix} \hat{\mathbf{B}}^b \\ \hat{\mathbf{u}} \end{pmatrix} - \begin{pmatrix} \mathbf{B}_k^b \\ \mathbf{u}_k \end{pmatrix} \right\|_2 \leq \frac{\varepsilon'}{2}.$$

Figure 2. Projected Stochastic Gradient Descent given access to samples from $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$.

- 1: $\mathbf{w}^{(0)} = \begin{pmatrix} (\mathbf{B}^{(0)})^b \\ \mathbf{u}^{(0)} \end{pmatrix} \leftarrow \begin{pmatrix} (\tilde{\boldsymbol{\Sigma}}_S^{-1})^b \\ \tilde{\boldsymbol{\mu}}_S \end{pmatrix}$
- 2: **for** $i = 1, \dots, T$ **do**
- 3: Sample $\mathbf{x}^{(i)}$ from $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*, S)$
- 4: $\eta_i \leftarrow \frac{1}{\lambda \cdot i}$
- 5: $\begin{pmatrix} (\mathbf{B}^{(i-1)})^b \\ \mathbf{u}^{(i-1)} \end{pmatrix} \leftarrow \mathbf{w}^{(i-1)}$
- 6:
$$\mathbf{v}^{(i)} \leftarrow \begin{pmatrix} \frac{1}{2} \left(\mathbf{x}^{(i)} \mathbf{x}^{(i)T} - \tilde{\boldsymbol{\Sigma}}_S - \tilde{\boldsymbol{\mu}}_S \tilde{\boldsymbol{\mu}}_S^T \right)^b \\ \tilde{\boldsymbol{\mu}}_S - \mathbf{x}^{(i)} \end{pmatrix}$$

$$e^{h(\mathbf{u}^{(i-1)}, \mathbf{B}^{(i-1)}; \mathbf{x}^{(i)})} \mathcal{N}(\mathbf{0}, \mathbf{I}; \mathbf{x}^{(i)}) \psi_k(\mathbf{x}^{(i)})$$
- 7: {From Equation (13).}
- 8: $\mathbf{r}^{(i)} \leftarrow \mathbf{w}^{(i-1)} - \eta_i \mathbf{v}^{(i)}$
- 9: $\mathbf{w}^{(i)} \leftarrow \arg \min_{\mathbf{w} \in \mathcal{D}} \|\mathbf{w} - \mathbf{r}^{(i)}\|_2^2$
 {From Lemma 8 of [DGTZ18].}
- 10: **end for**
- 11: $\begin{pmatrix} \hat{\mathbf{B}}^b \\ \hat{\mathbf{u}} \end{pmatrix} \leftarrow \frac{1}{T} \sum_{i=1}^T \mathbf{w}^{(i)}$
- 12: $\hat{\boldsymbol{\Sigma}} \leftarrow \hat{\mathbf{B}}^{-1}$
- 13: $\hat{\boldsymbol{\mu}} \leftarrow \hat{\mathbf{B}}^{-1} \hat{\mathbf{u}}$
- 14: **return** $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$

which together with (16) implies

$$\left\| \begin{pmatrix} \hat{\mathbf{B}}^b \\ \hat{\mathbf{u}} \end{pmatrix} - \begin{pmatrix} (\boldsymbol{\Sigma}^{*-1})^b \\ \boldsymbol{\mu}^* \end{pmatrix} \right\|_2 \leq \frac{\varepsilon'}{2}.$$

and the theorem follows as closeness in parameter distance implies closeness in total variation distance for the corresponding untruncated Gaussian distributions. \square

V. LOWER BOUND FOR LEARNING THE MEAN OF A TRUNCATED NORMAL

Theorem 7. *There exists a family of sets \mathcal{S} with $\Gamma(\mathcal{S}) = O(d)$ such that any algorithm that draws m samples from $\mathcal{N}(\boldsymbol{\mu}, \mathbf{I}, S)$ and computes an estimate $\tilde{\boldsymbol{\mu}}$ with $\|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2 \leq 1$ must have $m = \Omega(2^{d/2})$.*

Proof: Let $H = [-1, 1]^{d+1}$ be the $d + 1$ -dimensional cube. We will also use the left and right subcubes $H_+ = [-1, 0] \times [-1, 1]^d$, $H_- = [0, 1] \times [-1, 1]^d$ respectively. Let $\mathcal{N}_+ = \mathcal{N}(\mathbf{e}_1, \mathbf{I})$ and $\mathcal{N}_- = \mathcal{N}(-\mathbf{e}_1, \mathbf{I})$. We denote by r the (scaled) pointwise minimum of the two densities truncated

at the cube H , that is

$$\begin{aligned} r(\mathbf{x}) &= \frac{\min(\mathcal{N}_+(H; \mathbf{x}), \mathcal{N}_-(H; \mathbf{x}))}{c} \\ &= \frac{\mathbf{1}_H(\mathbf{x})}{c} \min(\mathcal{N}_+(\mathbf{x}), \mathcal{N}_-(\mathbf{x})), \end{aligned}$$

where $c = 1 - d_{\text{TV}}(\mathcal{N}_+, \mathcal{N}_-)$.

To simplify notation we assume that we work in \mathbb{R}^{d+1} instead of \mathbb{R}^d . Let $V = (v_1, \dots, v_d) \in \{+1, -1\}^d$. For every V we define the set $G_V = H \cap \{\mathbf{y} \in \mathbb{R}^d : y_i v_i \geq 0\}$. We also define the subcubes $H_V = [0, 1] \times G_V$. We consider the following subset of H parameterized by the 2^d parameters $t_V \in [0, 1]$ and $\delta \in [-1, 1]$.

$$S_+ = [-1 + \delta, 0] \times [-1, 1]^d \cup \bigcup_{V \in \{-1, +1\}^d} [0, t_V] \times G_V$$

We will argue that there exists a distribution D^+ on the values t_V such that on expectation $d_{\text{TV}}(\mathcal{N}_+^{S_+}, \mathcal{N}_-^{S_-})$ is $O(2^{-d})$. We show how to construct the distribution D_+ since the construction for D_- is the same. In fact we will show that both distributions are very close to $r(\mathbf{x})$. Notice that for some $(t, \mathbf{y}) \in \mathbb{R}^{d+1}$ we have We draw each t_V independently from the distribution with cdf

$$F(t) = \mathbf{1}_{[0,1)}(t)(1 - e^{-2t}) + \mathbf{1}_{[1,+\infty)}(t)$$

Notice that for $t \in (0, 1)$ and any $\mathbf{y} \in \mathbb{R}^d$ we have that $1 - F(t) = \mathcal{N}_-(t, \mathbf{y}) / \mathcal{N}_+(t, \mathbf{y})$.

After we draw all t_V from F we choose δ so that $\mathcal{N}_+(S_+; \mathbf{x}) = c$. We will show that on expectation over the t_V we have $\delta = 0$, which means that no correction is needed. In fact we show something stronger, namely that for all $\mathbf{x} \in H_+$ we have that $\mathbb{E}_{S_+ \sim D_+}[\mathcal{N}_+(S_+; \mathbf{x})] = r(\mathbf{x})$. Assume that $\mathbf{x} \in H_V$. Indeed,

$$\begin{aligned} &\mathbb{E}_{S_+ \sim D_+}[\mathcal{N}_+(S_+; \mathbf{x})] \\ &= \frac{\mathcal{N}_+(\mathbf{x})}{c} \mathbb{E}_{S_+ \sim D_+}[\mathbf{1}_{S_+}(\mathbf{x})] = \frac{\mathcal{N}_+(\mathbf{x})}{c} \mathbb{E}_{S_+ \sim D_+}[\mathbf{1}_{\{x_1 \leq t_V\}}] \\ &= \frac{\mathcal{N}_+(\mathbf{x})}{c} (1 - F(t_V)) = \frac{\mathcal{N}_-(\mathbf{x})}{c} = r(\mathbf{x}) \end{aligned}$$

Moreover, observe that for all $\mathbf{x} \in H_- \cap S_+$ we have that $\mathcal{N}_+(S_+; \mathbf{x}) = r(\mathbf{x})$ always (with probability 1). We now argue that in order to have constant probability to distinguish $\mathcal{N}_+(S_+)$ from $r(\mathbf{x})$ one needs to draw $\Omega(2^d)$ samples. Since the expected density of $\mathcal{N}_+(S_+)$ matches $r(\mathbf{x})$ for all $\mathbf{x} \in H_+$, to be able to distinguish the two distributions one needs to observe at least two samples in the same cube H_V . Since we have 2^d disjoint cubes H_V the probability of a sample landing in H_V is at most $1/2^d$. Therefore, using the birthday problem, to have constant probability to observe a collision one needs to draw $\Omega(\sqrt{2^d}) = \Omega(2^{d/2})$ samples. Since for all $\mathbf{x} \in H_- \cap S_+$, $\mathcal{N}_+(S_+)$ exactly matches $r(\mathbf{x})$, to distinguish between the two distributions one needs to observe a sample \mathbf{x} with $-1 + \delta < x_1 < -1$. Due to symmetry, \mathcal{N}_+ assigns

to all cubes H_V equal probability, call that p . Moreover, we have that $c = 2^{d+1}p$. Now let p_V be the random variable corresponding to the probability that N_+ assigns to $[0, t_V] \times G_V$. We have that $\mathbb{E}_{t_V \sim F}[p_V] = p$ for all V . Since the independent random variables p_V are bounded in $[0, 1/2^d]$, Hoeffding's inequality implies that $|\sum_{V \in \{-1, 1\}^d} (p_V - p)| < 1/2^{d/2}$ with probability at least $1 - 2/e^2$. This means that with probability at least $3/4$ one will need to draw $\Omega(2^{d/2})$ samples in order to observe one with $x_1 < -1 + \delta$.

Since any set S in our family \mathcal{S} has almost everywhere (that is except the set of its vertices which a finite set and thus of measure zero) smooth boundary we may use the following equivalent (see e.g. [Naz03]) definition of its surface area

$$\Gamma(S) = \int_{\partial S} \mathcal{N}_0(\mathbf{x}) d\sigma(\mathbf{x}),$$

where $d\sigma(\mathbf{x})$ is the standard surface measure on \mathbb{R}^d . Without loss of generality we assume that S corresponds to the set S_+ defined above (the proof is the same if we consider a set S_-). We have

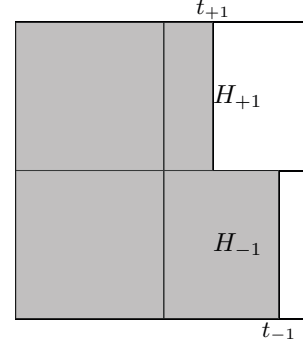
$$\begin{aligned} \partial S \subseteq & \bigcup_{V \in \{+1, -1\}^d} (\{t_V\} \times G_V) \\ & \cup \partial([-1, +1]^{d+1}) \cup \bigcup_{i=1}^{d+1} \{\mathbf{x} : x_i = 0\}. \end{aligned}$$

By the definition of Gaussian surface area it is clear that $\Gamma(A \cup B) \leq \Gamma(A) + \Gamma(B)$. From Table I we know that $\Gamma([-1, +1]^{d+1}) = O(\sqrt{\log d})$. Moreover, we know that a single halfspace has surface area at most $\sqrt{2/\pi}$ (see e.g. [KOS08]). Therefore $\Gamma(\bigcup_{i=1}^{d+1} \{\mathbf{x} : x_i = 0\}) \leq \sum_{i=1}^{d+1} \sqrt{2/\pi} = O(d)$. Finally, we notice that for any point \mathbf{x} on the hyperplane $\{\mathbf{x} : x_1 = 0\}$ and any \mathbf{y} on $\{\mathbf{x} : x_1 = c\}$ (for any $c \geq 0$), we have $\mathcal{N}_0(\mathbf{x}) \geq \mathcal{N}_0(\mathbf{y})$. Therefore, the surface area of each set $t_V \times G_V$ is maximized for $t_V = 0$. In this case $\bigcup_{V \in \{+1, -1\}^d} (\{t_V\} \times G_V) \subseteq \{\mathbf{x} : x_1 = 0\}$, which implies that the set $\bigcup_{V \in \{+1, -1\}^d} (\{t_V\} \times G_V)$ contributes at most $\sqrt{2/\pi}$ to the total surface area. Putting everything together, we have that $\Gamma(S) = O(d)$. ■

VI. IDENTIFIABILITY WITH BOUNDED GAUSSIAN SURFACE AREA

In this section we investigate the sample complexity of the problem of estimating the parameters of a truncated Gaussian using a different approach that does not depend on the VC dimension of the family \mathcal{S} of the truncation sets to be finite. For example, we settle the sample complexity of learning the parameters of a Gaussian distribution truncated at an unknown convex set (recall that the class of convex sets has infinite VC dimension). Our method relies on finding a tuple $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}, \tilde{S})$ of parameters so that the moments of

Figure 3. The set S_+ when $d = 1$.



the corresponding truncated Gaussian $\mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}, \tilde{S})$ are all close to the moments of the unknown truncated Gaussian distribution, for which we have unbiased estimates using samples. The main question that we need to answer to determine the sample complexity of this problem is how many moments are needed to be matched in order to be sure that our guessed parameters are close to the parameters of the unknown truncated Gaussian. We state now the main result. Its proof is based on Lemma 20 and can be found in the full version of the paper.

Theorem 8 (Moment Matching). *Let \mathcal{S} be a family of subsets of \mathbb{R}^d of bounded Gaussian surface area $\Gamma(\mathcal{S})$. Moreover, assume that if T is an affine map and $T(\mathcal{S}) = \{T(S) : S \in \mathcal{S}\}$ is the family of the images of the sets of \mathcal{S} , then it holds $\Gamma(T(\mathcal{S})) = O(\Gamma(\mathcal{S}))$. For some $S \in \mathcal{S}$, let $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ be an unknown truncated Gaussian. $d^{O(\Gamma(\mathcal{S})/\varepsilon^4)}$ samples are sufficient to find parameters $\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}, \tilde{S}$ such that $d_{\text{TV}}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S), \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}, \tilde{S})) \leq \varepsilon$.*

The key lemma of this section is Lemma 20. It shows that if two truncated normals are in total variation distance ε then there exists a moment where they differ. The main idea is to prove that there exists a polynomial that approximates well the indicator of the set $\{f_1 > f_2\}$. Notice that the total variation distance between two densities can be written as $\int \mathbf{1}_{\{f_1 > f_2\}}(\mathbf{x}) f_1(\mathbf{x}) - f_2(\mathbf{x}) d\mathbf{x}$. In our proof we use the chi squared divergence, which for two distributions with densities f_1, f_2 is defined as

$$D_{\chi^2}(f_1 \| f_2) = \int \frac{(f_1(\mathbf{x}) - f_2(\mathbf{x}))^2}{f_2(\mathbf{x})} d\mathbf{x}$$

To prove it we need the following nice fact about chi squared divergence between Gaussian distributions. In general chi squared divergence may be infinite for some pairs of Gaussians. In the following lemma we prove that for any pair of Gaussians, there exists another Gaussian N such that $D_{\chi^2}(N_1 \| N) D_{\chi^2}(N_2 \| N)$ are finite even if $D_{\chi^2}(N_1 \| N_2) = \infty$.

Lemma 19. *Let $N_1 = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$, and $N_2 = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_2)$ be two Normal distributions that satisfy the conditions of*

Lemma 3. Then there exists a Normal distribution N such that

$$D_{\chi^2}(N_1\|N), D_{\chi^2}(N_2\|N) \leq \exp\left(2\left\|\Sigma_1^{-1/2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\right\|_2\right) + \frac{1}{2} \max(1, \|\Sigma_1\|_2) \left\|\Sigma_1^{-1/2}\Sigma_2\Sigma_1^{-1/2} - \mathbf{I}\right\|_F^2$$

Now we state the main lemma of this section. We give here a sketch of its proof. For a formal proof we refer the reader to the full version of the paper.

Lemma 20. Let \mathcal{S} be a family of subsets of \mathbb{R}^d of bounded Gaussian surface area $\Gamma(\mathcal{S})$. Moreover, assume that if T is an affine map and $T(\mathcal{S}) = \{T(S) : S \in \mathcal{S}\}$ is the family of the images of the sets of \mathcal{S} , then it holds $\Gamma(T(\mathcal{S})) = O(\Gamma(\mathcal{S}))$. Let $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1, S_1)$ and $\mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2, S_2)$ be two truncated Gaussians with densities f_1, f_2 respectively. Let $k = O(\Gamma(\mathcal{S})/\varepsilon^4)$. If $d_{\text{TV}}(f_1, f_2) \geq \varepsilon$, then there exists a $V \in \mathbb{N}^d$ with $|V| \leq k$ such that

$$\left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1, S_1)}[\mathbf{x}^V] - \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2, S_2)}[\mathbf{x}^V] \right| \geq \varepsilon/d^{O(k)}.$$

Proof sketch. Let $W = S_1 \cap S_2 \cap \{f_1 > f_2\} \cup S_1 \setminus S_2$, that is the set of points where the first density is larger than the second. We now write the L_1 distance between f_1, f_2 as

$$\int |f_1(\mathbf{x}) - f_2(\mathbf{x})| d\mathbf{x} = \int \mathbf{1}_W(\mathbf{x})(f_1(\mathbf{x}) - f_2(\mathbf{x})) d\mathbf{x}$$

Denote $p(\mathbf{x})$ the polynomial that will do the approximation of the L_1 distance. From Lemma 19 we know that there exists a Normal distribution within small chi-squared divergence of both $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$ and $\mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2)$. Call the density function of this distribution $g(\mathbf{x})$. We have

$$\begin{aligned} & \left| \int |f_1(\mathbf{x}) - f_2(\mathbf{x})| d\mathbf{x} - \int p(\mathbf{x})(f_1(\mathbf{x}) - f_2(\mathbf{x})) d\mathbf{x} \right| \quad (18) \\ & \leq \int |\mathbf{1}_W(\mathbf{x}) - p(\mathbf{x})| |f_1(\mathbf{x}) - f_2(\mathbf{x})| d\mathbf{x} \\ & \leq \int |\mathbf{1}_W(\mathbf{x}) - p(\mathbf{x})| \sqrt{g(\mathbf{x})} \frac{|f_1(\mathbf{x}) - f_2(\mathbf{x})|}{\sqrt{g(\mathbf{x})}} d\mathbf{x} \\ & \leq \sqrt{\int (\mathbf{1}_W(\mathbf{x}) - p(\mathbf{x}))^2 g(\mathbf{x}) d\mathbf{x}} \\ & \cdot \sqrt{\int \frac{(f_1(\mathbf{x}) - f_2(\mathbf{x}))^2}{g(\mathbf{x})} d\mathbf{x}}, \quad (19) \end{aligned}$$

where we use Schwarz's inequality. From Lemma 19 we know that

$$\int \frac{f_1(\mathbf{x})^2}{g(\mathbf{x})} d\mathbf{x} \leq \int \frac{\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1; \mathbf{x})^2}{g(\mathbf{x})} d\mathbf{x} = \exp(\text{poly}(1/\alpha)).$$

Similarly, $\int \frac{f_2(\mathbf{x})^2}{g(\mathbf{x})} d\mathbf{x} = \exp(\text{poly}(1/\alpha))$. Therefore, we

have,

$$\begin{aligned} & \left| \int |f_1(\mathbf{x}) - f_2(\mathbf{x})| d\mathbf{x} - \int p(\mathbf{x})(f_1(\mathbf{x}) - f_2(\mathbf{x})) d\mathbf{x} \right| \\ & \leq \exp(\text{poly}(1/\alpha)) \sqrt{\int (\mathbf{1}_W(\mathbf{x}) - p(\mathbf{x}))^2 g(\mathbf{x}) d\mathbf{x}} \end{aligned}$$

Recall that $g(\mathbf{x})$ is the density function of a Gaussian distribution, and let $\boldsymbol{\mu}, \Sigma$ be the parameters of this Gaussian. Notice that it remains to show that there exists a good approximating polynomial $p(\mathbf{x})$ to the indicator function $\mathbf{1}_W$. We can now transform the space so that $g(\mathbf{x})$ becomes the standard normal. Notice that this is an affine transformation that also transforms the set W ; Since the Gaussian surface area is "invariant" under linear transformations

Since $\mathbf{1}_W \in L^2(\mathbb{R}^d, \mathcal{N}_0)$ we can approximate it using Hermite polynomials. For some $k \in \mathcal{N}$ we set $p(\mathbf{x}) = S_k \mathbf{1}_W(\mathbf{x})$, that is

$$p_k(\mathbf{x}) = \sum_{V: |V| \leq k} \widehat{\mathbf{1}_W} H_V(\mathbf{x}).$$

Combining Lemma 10 and Lemma 4 we obtain

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}_0} [(\mathbf{1}_W(\mathbf{x}) - p_k(\mathbf{x}))^2] = O\left(\frac{\Gamma(\mathcal{S})}{k^{1/2}}\right).$$

Therefore, $\left| \int |f_1(\mathbf{x}) - f_2(\mathbf{x})| d\mathbf{x} - \int p_k(\mathbf{x})(f_1(\mathbf{x}) - f_2(\mathbf{x})) d\mathbf{x} \right| = \exp(\text{poly}(1/\alpha)) \frac{\Gamma(\mathcal{S})^{1/2}}{k^{1/4}}$. Ignoring the dependence on the absolute constant α , to achieve error $O(\varepsilon)$ we need degree $k = O(\Gamma(\mathcal{S})^2/\varepsilon^4)$.

To complete the proof, it remains to obtain a bound for the coefficients of the polynomial $q(\mathbf{x}) = p_k(\Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}))$. Using known facts about the coefficients of Hermite polynomials we obtain that $\|q(\mathbf{x})\|_\infty \leq \binom{d+k}{k}^2 (4d)^{k/2} (O(1/\alpha^2))^k$. To conclude the proof we notice that we can pick the degree k so that

$$\begin{aligned} & \left| \int q(\mathbf{x})(f_1(\mathbf{x}) - f_2(\mathbf{x})) d\mathbf{x} \right| \\ & = \left| \sum_{V: |V| \leq k} \mathbf{x}^V (f_1(\mathbf{x}) - f_2(\mathbf{x})) \right| \geq \varepsilon/2. \end{aligned}$$

Since the maximum coefficient of $q(\mathbf{x})$ is bounded by $d^{O(k)}$ we obtain the result. \blacksquare

VII. VC-DIMENSION VS GAUSSIAN SURFACE AREA

We use two different complexity measures of the truncation set to get sample complexity bounds, the VC-dimension and the Gaussian Surface Area (GSA) of the class of the sets. As we already mentioned in the introduction, there are classes, for example convex sets, that have bounded Gaussian surface area but infinite VC-dimension. However, this is not the main difference between the two complexity

measures in our setting. Having a class with bounded VC-dimension means that the empirical risk minimization needs finite samples. To get an efficient algorithm we still need to *implement the ERM for this specific class*. Therefore, it is not clear whether it is possible to get an algorithm that works for all sets of bounded VC-dimension. On the other hand, bounded GSA means that we can approximate the weighted indicator of the set using its low order Hermite coefficients. This approximation works for all sets of bounded GSA and does not depend on the specific class of sets. Therefore, using GSA we manage to get a unified approach that learns the parameters of the underlying Gaussian distribution using only the assumption that the truncation set has bounded GSA. In other words, our approach uses the information of the class that the truncation set belongs only to decide how large the degree of the approximating polynomial should be. Having said that, it is an interesting open problem to design algorithms that learn the parameters of the Gaussian and use the information that the truncation set belongs to some class (e.g. intersection of k -halfspaces) to beat the runtime of our generic approach that only depends on the GSA of the class.

REFERENCES

- [AGR13] Joseph Anderson, Navin Goyal, and Luis Rademacher. Efficient learning of simplices. In *Conference on Learning Theory*, pages 1020–1045, 2013.
- [Bal93] Keith Ball. The reverse isoperimetric problem for gaussian measure. *Discrete & Computational Geometry*, 10(1):411–420, 1993.
- [BC14] N Balakrishnan and Erhard Cramer. *The art of progressive censoring*. Springer, 2014.
- [Coh16] A Clifford Cohen. *Truncated and censored samples: theory and applications*. CRC press, 2016.
- [CSV17] Moses Charikar, Jacob Steinhardt, and Gregory Valiant. Learning from untrusted data. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017, Montreal, QC, Canada, June 19-23, 2017*, pages 47–60, 2017.
- [CW01] Anthony Carbery and James Wright. Distributional and l_q norm inequalities for polynomials over convex bodies in \mathbb{R}^n . *Mathematical research letters*, 8(3):233–248, 2001.
- [DDS14] Anindya De, Ilias Diakonikolas, and Rocco A Servedio. Learning from satisfying assignments. In *Proceedings of the twenty-sixth annual ACM-SIAM symposium on Discrete algorithms*, pages 478–497. SIAM, 2014.
- [Den98] François Denis. Pac learning from positive statistical queries. In *International Conference on Algorithmic Learning Theory*, pages 112–126. Springer, 1998.
- [DGTZ18] Constantinos Daskalakis, Themis Gouleakis, Christos Tzamos, and Manolis Zampetakis. Efficient statistics, in high dimensions, from truncated samples. In *the 59th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 2018.
- [DK14] Constantinos Daskalakis and Gautam Kamath. Faster and sample near-optimal algorithms for proper learning mixtures of gaussians. In *Proceedings of The 27th Conference on Learning Theory, COLT 2014, Barcelona, Spain, June 13-15, 2014*, pages 1183–1213, 2014.
- [DKK⁺16] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high dimensions without the computational intractability. In *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016, 9-11 October 2016, Hyatt Regency, New Brunswick, New Jersey, USA*, pages 655–664, 2016.
- [DKK⁺17] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Being robust (in high dimensions) can be practical. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 999–1008, 2017.
- [DKK⁺18] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robustly learning a gaussian: Getting optimal error, efficiently. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018, New Orleans, LA, USA, January 7-10, 2018*, pages 2683–2702, 2018.
- [DL12] Luc Devroye and Gábor Lugosi. *Combinatorial methods in density estimation*. Springer Science & Business Media, 2012.
- [Eld11] Ronen Eldan. A polynomial number of random points does not determine the volume of a convex body. *Discrete & Computational Geometry*, 46(1):29–47, 2011.
- [Fis31] RA Fisher. Properties and applications of Hh functions. *Mathematical tables*, 1:815–852, 1931.
- [FJK96] Alan Frieze, Mark Jerrum, and Ravi Kannan. Learning linear transformations. In *Foundations of Computer Science, 1996. Proceedings., 37th Annual Symposium on*, pages 359–368. IEEE, 1996.
- [Gal97] Francis Galton. An examination into the registered speeds of american trotting horses, with remarks on their value as hereditary data. *Proceedings of the Royal Society of London*, 62(379-387):310–315, 1897.
- [GR09] Navin Goyal and Luis Rademacher. Learning convex bodies is hard. *arXiv preprint arXiv:0904.1227*, 2009.
- [Kan11] Daniel M Kane. The gaussian surface area and noise sensitivity of degree-d polynomial threshold functions. *computational complexity*, 20(2):389–412, 2011.
- [KKMS05] Adam Tauman Kalai, Adam R. Klivans, Yishay Mansour, and Rocco A. Servedio. Agnostically learning halfspaces. In *46th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2005), 23-25 October 2005, Pittsburgh, PA, USA, Proceedings*, pages 11–20, 2005.

- [KOS08] Adam R. Klivans, Ryan O’Donnell, and Rocco A. Servedio. Learning geometric concepts via gaussian surface area. In *49th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2008, October 25-28, 2008, Philadelphia, PA, USA*, pages 541–550, 2008.
- [LDG00] Fabien Letouzey, François Denis, and Rémi Gilleron. Learning from positive and unlabeled examples. In *International Conference on Algorithmic Learning Theory*, pages 71–85. Springer, 2000.
- [Led94] Michel Ledoux. Semigroup proofs of the isoperimetric inequality in euclidean and gauss space. *Bulletin des sciences mathématiques*, 118(6):485–510, 1994.
- [Lee14] Alice Lee. Table of the gaussian” tail” functions; when the” tail” is larger than the body. *Biometrika*, 10(2/3):208–214, 1914.
- [LRV16] Kevin A. Lai, Anup B. Rao, and Santosh Vempala. Agnostic estimation of mean and covariance. In *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016, 9-11 October 2016, Hyatt Regency, New Brunswick, New Jersey, USA*, pages 665–674, 2016.
- [Naz03] Fedor Nazarov. *On the Maximal Perimeter of a Convex Set in \mathbb{R}^n with Respect to a Gaussian Measure*, pages 169–187. Springer Berlin Heidelberg, Berlin, Heidelberg, 2003.
- [O’D14] Ryan O’Donnell. *Analysis of Boolean Functions*. Cambridge University Press, 2014.
- [Pea02] Karl Pearson. On the systematic fitting of frequency curves. *Biometrika*, 2:2–7, 1902.
- [Pis86] Gilles Pisier. Probabilistic methods in the geometry of banach spaces. In *Probability and analysis*, pages 167–241. Springer, 1986.
- [PL08] Karl Pearson and Alice Lee. On the generalised probable error in multiple normal correlation. *Biometrika*, 6(1):59–68, 1908.
- [Sch86] Helmut Schneider. *Truncated and censored samples from normal populations*. Marcel Dekker, Inc., 1986.
- [SJ66] SM Shah and MC Jaiswal. Estimation of parameters of doubly truncated normal distribution from first four sample moments. *Annals of the Institute of Statistical Mathematics*, 18(1):107–111, 1966.
- [SSBD14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [Sze67] G. Szegö. *Orthogonal Polynomials*. Number τ . 23 in American Mathematical Society colloquium publications. American Mathematical Society, 1967.