

# Robustly Learning Mixtures of $k$ Arbitrary Gaussians

Ainesh Bakshi  
CMU  
abakshi@cs.cmu.edu

Ilias Diakonikolas  
UW Madison  
ilias@cs.wisc.edu

He Jia  
Georgia Tech  
hjia36@gatech.edu

Daniel M. Kane  
UCSD  
dakane@ucsd.edu

Pravesh K. Kothari  
CMU  
praveshk@cs.cmu.edu

Santosh S. Vempala  
Georgia Tech  
vempala@gatech.edu

June 8, 2021

## Abstract

We give a polynomial-time algorithm for the problem of robustly estimating a mixture of  $k$  arbitrary Gaussians in  $\mathbb{R}^d$ , for any fixed  $k$ , in the presence of a constant fraction of arbitrary corruptions. This resolves the main open problem in several previous works on algorithmic robust statistics, which addressed the special cases of robustly estimating (a) a single Gaussian, (b) a mixture of TV-distance separated Gaussians, and (c) a uniform mixture of two Gaussians. Our main tools are an efficient *partial clustering* algorithm that relies on the sum-of-squares method, and a novel *tensor decomposition* algorithm that allows errors in both Frobenius norm and low-rank terms.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background and Motivation . . . . .	1
1.2	Our Results . . . . .	2
1.3	Concurrent and Independent Work . . . . .	6
1.4	Related Prior Work . . . . .	7
1.5	Organization . . . . .	8
1.6	Overview of Techniques . . . . .	8
<b>2</b>	<b>Preliminaries</b>	<b>16</b>
2.1	Gaussian Background . . . . .	16
2.2	Sum-of-Squares Proofs and Pseudo-distributions . . . . .	18
2.3	Analytic Properties of Gaussian Distributions . . . . .	22
2.4	Deterministic Conditions on the Uncorrupted Samples . . . . .	26
2.5	Hypothesis Selection . . . . .	28
<b>3</b>	<b>List-Recovery of Parameters via Tensor Decomposition</b>	<b>28</b>
3.1	List-Decodable Tensor Decomposition Algorithm . . . . .	28
3.2	Analysis of Algorithm . . . . .	30
3.3	Robust Estimation of Hermite Tensors . . . . .	31
3.4	List-Recovery of Covariances up to Low-Rank Error . . . . .	32
3.5	Finding a Low-dimensional Subspace for Exhaustive Search . . . . .	39
3.6	Parameter vs Moment Distance for Gaussian Mixtures . . . . .	42
<b>4</b>	<b>Robust Partial Cluster Recovery</b>	<b>44</b>
4.1	Algorithm . . . . .	46
4.2	Analysis . . . . .	47
4.3	Proof of Lemma 4.5 . . . . .	50
4.4	Proof of Lemma 4.6 . . . . .	51
4.5	Special Case: Algorithm for Uniform and Bounded Mixing Weights . . . . .	54
<b>5</b>	<b>Spectral Separation of Thin Components</b>	<b>55</b>
<b>6</b>	<b>Robust Proper Learning: Proof of Theorem 1.4</b>	<b>58</b>
6.1	Analysis of Algorithm 6.3 . . . . .	61
6.2	Proof of the Main Theorem . . . . .	65
<b>7</b>	<b>More Efficient Robust Partial Cluster Recovery</b>	<b>69</b>
7.1	Algorithm . . . . .	70
7.2	Proof of Lemma 7.4 . . . . .	75
7.3	2nd Moment Estimation Subroutine . . . . .	78
<b>8</b>	<b>Getting <math>\text{poly}(\epsilon)</math>-close in TV Distance: Proof of Theorem 1.5</b>	<b>80</b>
8.1	Proof of Theorem 8.2 . . . . .	81
8.2	Proof of Theorem 8.1 . . . . .	86
<b>9</b>	<b>Robust Parameter Recovery: Proof of Theorem 1.6</b>	<b>88</b>
9.1	Proof of Lemma 9.9 . . . . .	96

<b>References</b>	<b>99</b>
<b>A Omitted Proofs</b>	<b>104</b>
A.1 Omitted Proofs from Section 2.1 . . . . .	104
A.2 Omitted Proofs from Section 2.2 . . . . .	105
A.3 Omitted Proofs from Section 2.3 . . . . .	105
A.4 Omitted Proofs from Section 2.4 . . . . .	111
A.5 Omitted Proofs from Section 6 . . . . .	113
<b>B Bit Complexity Analysis</b>	<b>114</b>

# 1 Introduction

## 1.1 Background and Motivation

Given a collection of observations and a class of models, the objective of a typical learning algorithm is to find the model in the class that best fits the data. The classical assumption is that the input data are i.i.d. samples generated by a statistical model in the given class. This is a simplifying assumption that is, at best, only approximately valid, as real datasets are typically exposed to some source of systematic noise. Robust statistics [HRRS86, HR09] challenges this assumption by focusing on the design of *outlier-robust* estimators — algorithms that can tolerate a *constant fraction* of corrupted datapoints, independent of the dimension. Despite significant effort over several decades starting with important early works of Tukey and Huber in the 60s, even for the most basic high-dimensional estimation tasks, all known computationally efficient estimators were until fairly recently highly sensitive to outliers.

This state of affairs changed with two independent works from the TCS community [DKK<sup>+</sup>16, LRV16], which gave the first computationally efficient and outlier-robust learning algorithms for a range of “simple” high-dimensional probabilistic models. In particular, these works developed efficient robust estimators for a single high-dimensional Gaussian distribution with unknown mean and covariance. Since these initial algorithmic works [DKK<sup>+</sup>16, LRV16], we have witnessed substantial research progress on algorithmic aspects of robust high-dimensional estimation by several communities of researchers, including TCS, machine learning, and mathematical statistics. See Section 1.4 for an overview of the prior work most relevant to the results of this paper. The reader is referred to [DK19] for a recent survey on the topic.

One of the main original motivations for the development of algorithmic robust statistics within the TCS community was the problem of learning high-dimensional Gaussian mixture models. A *Gaussian mixture model (GMM)* is a convex combination of Gaussian distributions, i.e., a distribution on  $\mathbb{R}^d$  of the form  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$ , where the weights  $w_i$ , mean vectors  $\mu_i$ , and covariance matrices  $\Sigma_i$  are unknown. GMMs are *the* most extensively studied latent variable model in the statistics and machine learning literatures, starting with the pioneering work of Karl Pearson in 1894 [Pea94], which introduced the method of moments in this context.

In the absence of outliers, a long line of work initiated by Dasgupta [Das99, AK01, VW02, AM05, KSV08, BV08] gave efficient clustering algorithms for GMMs under various separation assumptions. Subsequently, efficient learning algorithms were obtained [KMV10, MV10, BS10, HP15] under minimal information-theoretic conditions. Specifically, Moitra and Valiant [MV10] and Belkin and Sinha [BS10] designed the first polynomial-time learning algorithms for arbitrary Gaussian mixtures with any fixed number of components. These works qualitatively characterized the complexity of this fundamental learning problem in the noiseless setting. Alas, all aforementioned algorithms are very fragile in the presence of corrupted data. Specifically, a *single* outlier can completely compromise their performance.

Developing efficient learning algorithms for high-dimensional GMMs in the more realistic *outlier-robust* setting — the focus of the current paper — has turned out to be significantly more challenging. This was both one of the original motivations and the main open problem in the

initial robust statistics works [DKK<sup>+</sup>16, LRV16]. We note that [DKK<sup>+</sup>16] developed a robust density estimation algorithm for mixtures of *spherical* Gaussians — a very special case of our problem where the covariance of each component is a multiple of the identity — and highlighted a number of key technical obstacles that need to be overcome in order to handle the general case. Since then, a number of works have made algorithmic progress on important special cases of the general problem. These include faster robust clustering for the spherical case under minimal separation conditions [HL18, KSS18, DKS18], robust clustering for separated (and potentially non-spherical) Gaussian mixtures [BK20b, DHKK20], and robustly learning *uniform* mixtures of two arbitrary Gaussian components [Kan20].

This progress notwithstanding, the algorithmic task of robustly learning a mixture of a constant number (or even two) arbitrary Gaussians (with arbitrary weights) has remained a central open problem in this field, as highlighted recently [DVW19].

This discussion motivates the following question, whose resolution is the main result of this work:

**Question 1.1.** *Is there a  $\text{poly}(d, 1/\varepsilon)$ -time robust GMM learning algorithm, in the presence of an  $\varepsilon$ -fraction of outliers, that has a dimension-independent error guarantee, for an arbitrary mixture of any constant number of arbitrary Gaussians on  $\mathbb{R}^d$ ?*

## 1.2 Our Results

To formally state our main result, we define the model of robustness we study. We focus on the following standard data corruption model that generalizes Huber’s contamination model [Hub64].

**Definition 1.2** (Total Variation Contamination Model). Given a parameter  $0 < \varepsilon < 1/2$  and a class of distributions  $\mathcal{F}$  on  $\mathbb{R}^d$ , the *adversary* operates as follows: The algorithm specifies the number of samples  $n$ . The adversary knows the true target distribution  $X \in \mathcal{F}$  and selects a distribution  $F$  such that  $d_{\text{TV}}(F, X) \leq \varepsilon$ . Then  $n$  i.i.d. samples are drawn from  $F$  and are given as input to the algorithm.

Intuitively, the parameter  $\varepsilon$  in Definition 1.2 quantifies the power of the adversary. The total variation contamination model is strictly stronger than Huber’s contamination model. Recall that in Huber’s model [Hub64], the adversary generates samples from a mixture distribution  $F$  of the form  $F = (1 - \varepsilon)X + \varepsilon N$ , where  $X$  is the unknown target distribution and  $N$  is an adversarially chosen noise distribution. That is, in Huber’s model the adversary is only allowed to add outliers.

*Remark 1.3.* The *strong contamination model* [DKK<sup>+</sup>16] is a strengthening of the total variation contamination, where an adversary can see the clean samples and then arbitrarily replace an  $\varepsilon$ -fraction of these points to obtain an  $\varepsilon$ -corrupted set of samples. Our robust learning algorithm succeeds in this strong contamination model, with the additional requirement that we can obtain two sets of independent  $\varepsilon$ -corrupted samples from the unknown mixture.

In the context of robustly learning GMMs, we want to design an efficient algorithm with the following performance: Given a sufficiently large set of samples from a distribution that is  $\varepsilon$ -close in total variation distance to an unknown GMM  $\mathcal{M}$  on  $\mathbb{R}^d$ , the algorithm outputs a hypothesis GMM

$\widehat{\mathcal{M}}$  such that with high probability the total variation distance  $d_{\text{TV}}(\widehat{\mathcal{M}}, \mathcal{M})$  is small. Specifically, we want  $d_{\text{TV}}(\widehat{\mathcal{M}}, \mathcal{M})$  to be only a function of  $\varepsilon$  and independent of the underlying dimension  $d$ .

The main result of this paper is the following:

**Theorem 1.4** (Main Result, See Corollary 6.2). *There is an algorithm with the following behavior: Given  $\varepsilon > 0$  and a multiset of  $n = d^{O(k)} \text{poly}(\log(1/\varepsilon))$  samples from a distribution  $F$  on  $\mathbb{R}^d$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , for an unknown target  $k$ -GMM  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$ , the algorithm runs in time  $\text{poly}(n) \text{poly}_k(1/\varepsilon)$  and outputs a  $k$ -GMM hypothesis  $\widehat{\mathcal{M}} = \sum_{i=1}^k \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  such that with high probability we have that  $d_{\text{TV}}(\widehat{\mathcal{M}}, \mathcal{M}) \leq g(\varepsilon, k)$ . Here  $g : \mathbb{R}_+ \times \mathbb{Z}_+ \rightarrow \mathbb{R}_+$  is a function such that  $\lim_{\varepsilon \rightarrow 0} g(\varepsilon, k) = 0$ .*

Theorem 1.4 gives the first polynomial-time *robust proper learning* algorithm, with dimension-independent error guarantee, for *arbitrary*  $k$ -GMMs, for any fixed  $k$ . This is the first polynomial-time algorithm for this problem, even for  $k = 2$ .

Since the dissemination of [BDJ<sup>+</sup>20], we have improved the guarantees of Theorem 1.4 in two ways. First, by refining the guarantees of two of our components in the algorithm, we are able to obtain a robust proper learner with error guarantee of  $\text{poly}_k(\varepsilon)$  (Theorem 1.5). Second, we show that the same algorithm also achieves the stronger parameter estimation guarantee (Theorem 1.6). These refinements build heavily on the techniques in [BDJ<sup>+</sup>20] and represent work subsequent to [BDJ<sup>+</sup>20, LM20]. We describe them in detail in the following.

**Independent and Concurrent Work** Theorem 1.4 is the main result obtained in [BDJ<sup>+</sup>20]. In independent and concurrent work to [BDJ<sup>+</sup>20], Liu and Moitra [LM20] obtained a closely related result under stronger assumptions on the input mixture — specifically,  $1/f(k)$  lower bound on the component weights,  $\text{poly}(\varepsilon/d)$  lower bound and  $\text{poly}(d/\varepsilon)$  upper bound on the eigenvalues of each component covariance, and a  $\geq g(k)$  TV-distance separation between every pair of components. Under these assumptions, they obtain a formally stronger robust *parameter estimation* guarantee in time that grows exponentially in  $f(k), g(k)$  with error guarantees that decay exponentially in  $f(k), g(k)$ . We discuss this work and its connection to Theorem 1.4 in more detail in Section 1.3.

**Discussion** Before proceeding, we make a few important remarks about Theorem 1.4.

1. *Sample Complexity and Runtime:* Our algorithm succeeds whenever the sample size  $n$  satisfies  $n \geq n_0 = d^{O(k)}/\text{poly}(\varepsilon)$ . The running time of our algorithm is  $\text{poly}(n) \text{poly}_k(1/\varepsilon)$ . Statistical query lower bounds [DKS17] suggest that  $d^{\Omega(k)}$  samples are necessary for efficiently learning GMMs, even for approximation to constant accuracy in the simpler setting without outliers and under the more restrictive *clustering* setting (where components are pairwise well-separated in total variation distance). This provides some evidence that the sample-time tradeoff achieved by Theorem 1.4 is qualitatively optimal (within absolute constant factors in the exponent). We note that the algorithm establishing Theorem 1.4 works in the standard bit-complexity model of computation and its running time is polynomial in the bit-complexity of the input parameters. We discuss the numerical accuracy required to implement our algorithm in Appendix B.

In the noiseless case, the first polynomial-time learning algorithm for  $k$ -GMMs on  $\mathbb{R}^d$  was given in [MV10, BS10]. In particular, the sample complexity and running time of the [MV10] algorithm is  $(d/\varepsilon)^{q(k)}$ , for some function  $q(k) = k^{\Omega(k)}$ . We observe that our running time and sample complexity are exponentially better than the guarantees for the noiseless case in [MV10, BS10]. Moreover, as we explain in Section 1.6, the [MV10, BS10] algorithms are very sensitive to outliers and an entirely new approach is required to obtain an efficient robust learning algorithm.

2. *Handling Arbitrary Weights:* The algorithm of Theorem 1.4 succeeds *without any assumptions* on the weights of the mixture components. We emphasize that this is an important feature and not a technicality. Prior work [BK20b, DHKK20, Kan20], as well as the concurrent work [LM20], cannot handle the case of general weights — even for the case of  $k = 2$  components. In fact, for the special case of uniform weights, we give a simpler algorithm for robustly learning GMMs (presented in Theorem 4.12). This algorithm naturally generalizes to give a sample complexity and running time that grows *exponentially* in  $1/w_{\min}$ , where  $w_{\min}$  is the minimum weight of any component in the mixture. Handling the general case (i.e., obtaining a fully polynomial-time algorithm, not incurring an exponential cost in  $1/w_{\min}$ ) requires genuinely new algorithmic ideas and is one of the key technical innovations in the proof of Theorem 1.4.
3. *Handling Arbitrary Covariances:* The algorithm of Theorem 1.4 does not require assumptions on the variances of the component covariances, modulo basic limitations posed by numerical computation issues. Specifically, our algorithm works even if some of the component covariances are rank-deficient (i.e., have directions of 0 variance) with running time scaling polynomially in the bit-complexity of the unknown component means and covariances. Such a dependence on the bit complexity of the input parameters is unavoidable – there exist<sup>1</sup> examples of rank-deficient covariances with irrational entries such that the total variation distance between the corresponding Gaussian and every Gaussian with covariance matrix of rational entries is the maximum possible value of one.
4. *Error Guarantee:* The function  $g$  quantifying the final error guarantee of our basic algorithm is  $g(\varepsilon, k) = 1/(\log(1/\varepsilon))^{C_k}$ , for some function  $C_k$  that goes to 0 when  $k$  increases. Importantly, for any fixed  $k$ , the final error guarantee of our algorithm depends only on  $\varepsilon$ , tends to 0 as  $\varepsilon \rightarrow 0$  and is independent of the dimension  $d$ . In Theorem 1.5, we show that, by modifying our algorithm, we can obtain improved error – scaling as a fixed polynomial in  $\varepsilon$ . This turns out to be quantitatively close to best possible for any robust proper learning algorithm.

Our work is most closely related to the recent paper by Kane [Kan20], which gave a polynomial-time robust learning algorithm for the *uniform*  $k = 2$  case, i.e., the case of two *equal weight* components, and the polynomial time algorithms [BK20b, DHKK20] for the problem under the (strong) assumption that the component Gaussians are pairwise well-separated in total variation distance.

---

<sup>1</sup>For e.g., for unit vector  $v = (1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3}, 0, 0, \dots, 0)$  and for every choice of rational covariance  $\Sigma$ , the total variation distance between  $\mathcal{N}(0, I - vv^\top)$  and  $\mathcal{N}(0, \Sigma)$  is the maximum possible 1.

Our algorithm builds on the ideas in the works [BK20b, DHKK20] that gave efficient clustering algorithms for any fixed number  $k$  of components, under the crucial assumption that the components have pairwise total variation distance close to 1. In this case, the above works actually succeed in efficiently *clustering* the input sample into  $k$  groups, such that each group contains the samples generated from one of the Gaussians, up to some small misclassification error. In contrast, the main challenge in this work is the information-theoretic impossibility of clustering in our setting where there are no separation assumptions. As we will explain in the proceeding discussion, while we draw ideas from [BK20b, DHKK20, Kan20], a number of significant conceptual and technical challenges need to be overcome in the non-clusterable setting.

**Improvements to Theorem 1.4.** We now describe refinements of our main theorem.

**Improving Error to a Fixed Polynomial in  $\varepsilon$**  It turns out that the inverse poly-logarithmic accuracy (in  $1/\varepsilon$ ) in the final error guarantee of Theorem 1.4 can be traced to an exhaustive search subroutine in our novel tensor decomposition subroutine and probability of success of our rounding algorithm in our partial clustering routine. Via natural (and conceptually simple) quantitative improvements to these two ingredients, we obtain an algorithm achieving the qualitatively nearly best possible error of  $\text{poly}_k(\varepsilon)$ . Specifically, we show:

**Theorem 1.5** (Robustly Learning  $k$ -Mixtures with  $\text{poly}(\varepsilon)$ -error, Informal, see Corollary 8.1). *There is an algorithm with the following behavior: Given  $\varepsilon > 0$  and a multiset of  $n = d^{O(k)} \text{poly}_k(1/\varepsilon)$  samples from a distribution  $F$  on  $\mathbb{R}^d$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , for an unknown target  $k$ -GMM  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$ , the algorithm runs in time  $\text{poly}(n) \text{poly}_k(1/\varepsilon)$  and outputs a  $k$ -GMM hypothesis  $\widehat{\mathcal{M}} = \sum_{i=1}^k \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  such that with high probability we have that  $d_{\text{TV}}(\widehat{\mathcal{M}}, \mathcal{M}) \leq \mathcal{O}(\varepsilon^{c_k})$ , where  $c_k$  depends only on  $k$ .*

**Robust Parameter Recovery** Finally, we show that the *same* algorithm as in Theorems 1.4 and 8.1 actually implies that the recovered mixture of Gaussians is close in *parameter* distance to the unknown target mixture. Such parameter estimation results are usually stated under the assumption that every pair of components of the unknown mixture are separated in total variation distance. In this work, we provide a stronger version of this parameter estimation guarantee.

More specifically, in the theorem below, we prove that whenever the components of the input mixture can be clustered together into some groups such that all mixtures in a group are close (and thus, indistinguishable), there exists a similar clustering of the output mixture such that all parameters (weight, mean, and covariances) of each cluster are close within  $\text{poly}_k(\varepsilon)$  in total variation distance. In particular this means that for each significant component of the input mixture, there is a component of the output mixture with very close parameters.

We note that [LM20] gave a parameter estimation guarantee (under additional assumptions on the mixture weights and component variances) whenever every pair of components in the unknown mixture are  $f(k)$ -far in total variation distance, where  $f$  can be any function of  $k$ , but the choice of  $f$  affects the exponent in the running time and error guarantee of the [LM20] algorithm.)

By strengthening one of the structural results in their argument, we establish the following:



**Theorem 1.6** (Parameter Recovery, See Theorem 9.2). *Given  $\varepsilon > 0$  and a multiset of  $n = d^{O(k)} \text{poly}_k(1/\varepsilon)$  samples from a distribution  $F$  on  $\mathbb{R}^d$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , for an unknown target  $k$ -GMM  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$ , the algorithm runs in time  $\text{poly}(n) \text{poly}_k(1/\varepsilon)$  and outputs a  $k'$ -GMM hypothesis  $\widehat{\mathcal{M}} = \sum_{i=1}^{k'} \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  with  $k' \leq k$  such that with high probability we have that there exists a partition of  $[k]$  into  $k' + 1$  sets  $R_0, R_1, \dots, R_{k'}$  such that*

1. Let  $W_i = \sum_{j \in R_i} w_j$ ,  $i \in \{0, 1, \dots, k'\}$ . Then, for all  $i \in [k']$ , we have that

$$|W_i - \widehat{w}_i| \leq \text{poly}_k(\varepsilon), \text{ and}$$

$$d_{\text{TV}}(\mathcal{N}(\mu_j, \Sigma_j), \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)) \leq \text{poly}_k(\varepsilon) \quad \forall j \in R_i .$$

2. The total weight of exceptional components in  $R_0$  is  $W_0 \leq \text{poly}_k(\varepsilon)$ .

If we assume additionally that any pair of components in the unknown mixture has total variation distance at least  $\text{poly}_k(\varepsilon)$ , then the following result follows directly from Theorem 1.6.

**Corollary 1.7.** *Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  be an unknown target  $k$ -GMM satisfying the following conditions: (i)  $d_{\text{TV}}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j)) \geq \varepsilon^{f_1(k)}$  for all  $i \neq j$ , and (ii)  $S = \{i \in [k] : w_i \geq \varepsilon^{f_2(k)}\}$  is a subset of  $[k]$ , where  $f_1(k), f_2(k)$  are sufficiently small functions of  $k$ . Given  $\varepsilon > 0$  and a multiset of  $n = d^{O(k)} \text{poly}_k(1/\varepsilon)$  samples from a distribution  $F$  on  $\mathbb{R}^d$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , there exists an algorithm that runs in time  $\text{poly}(n) \text{poly}_k(1/\varepsilon)$  and outputs a  $k'$ -GMM hypothesis  $\widehat{\mathcal{M}} = \sum_{i=1}^{k'} \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  with  $k' \leq k$  such that with high probability there exists a bijection  $\pi : S \rightarrow [k']$  satisfying the following: For all  $i \in S$ , it holds that*

$$|w_i - \widehat{w}_{\pi(i)}| \leq \text{poly}_k(\varepsilon)$$

$$d_{\text{TV}}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\widehat{\mu}_{\pi(i)}, \widehat{\Sigma}_{\pi(i)})) \leq \text{poly}_k(\varepsilon).$$

We note that both the pairwise separation between the components and the lower bounds on the weights in Corollary 1.7 scale as a fixed polynomial in  $\varepsilon$  (for fixed  $k$ ), which is qualitatively information-theoretically necessary.

### 1.3 Concurrent and Independent Work

In independent and concurrent work with [BDJ<sup>+</sup>20], which established Theorem 1.4, [LM20] gave an efficient outlier-robust parameter learning algorithm for mixtures of Gaussians under additional assumptions. Unlike our main theorem (Theorem 1.4), the algorithm in [LM20] succeeds under three crucial assumptions: 1) all mixing weights are bounded below by  $1/f(k)$ , where  $f$  is a function of  $k$  that appears in the exponent of the running time and output error, 2) all components have covariances with eigenvalues at least  $\text{poly}(\varepsilon/d)$  and at most  $\text{poly}(d/\varepsilon)$  (in particular, they cannot handle rank-deficient components), and 3) every pair of component Gaussians are separated by  $g(k)$  in total variation distance (regardless of how small the fraction of corruptions  $\varepsilon$  is). (We again note that each of these obtaining an algorithm without these assumptions requires significant new ideas in both the design and the analysis of the algorithm. Our techniques lead to efficient

algorithms without any of these assumptions.) The analysis of [LM20] obtained a robust *parameter recovery* guarantee on the output of their algorithm. In contrast, our Theorem 1.4 does not make any assumption on the mixture weights, the covariances or component separations, yielding the weaker guarantee of *proper learning* (i.e., returning a mixture of Gaussians that is close in total variation distance to the unknown uncorrupted mixture). The error guarantee offered in our Theorem 1.4, while being a constant (for any fixed  $k, \varepsilon$ ) is inverse poly-logarithmic in  $1/\varepsilon$ . We subsequently were able to improve these guarantees to achieve  $\text{poly}_k(\varepsilon)$  error and parameter recovery guarantees.

We also note that the techniques adopted in the two concurrent works ([BDJ<sup>+</sup>20] and [LM20]) are significantly different. On the one hand, both works make essential use of recent advances in clustering non-spherical mixtures [BK20b, DHKK20]. While [LM20] relies on the clustering algorithm of [DHKK20], our work leverages a *key modification* to the clustering algorithm of [BK20b] that allows us to obtain a fixed polynomial-time algorithm (as opposed to incurring an exponential cost in  $1/w_{\min}$ , where  $w_{\min}$  is the minimum mixing weight), at the cost of handling only a *subclass* of clusterable Gaussian mixtures. Indeed, our novel variant of clustering combined with our new spectral clustering subroutine is the key to not incurring an exponential cost in  $1/w_{\min}$ . In the next step, where the input mixture can be assumed to be non-clusterable, [LM20] uses a new sum-of-squares based algorithm for parameter estimation. In contrast, our work does not use sum-of-squares method in this component and instead relies on a new list-decoding algorithm for tensor decomposition that makes no assumption on the underlying components.

## 1.4 Related Prior Work

The algorithmic question of designing efficient robust estimators in high dimensions has been extensively studied in recent years. After the initial papers [DKK<sup>+</sup>16, LRV16], a number of works developed robust estimators for a range of statistical problems. These include efficient outlier-robust algorithms for sparse estimation [BDLS17, DKK<sup>+</sup>19], learning graphical models [CDKS18], linear regression [KKM18, DKS19, BP20, ZJS20, CAT<sup>+</sup>20], stochastic optimization [PSBR18, DKK<sup>+</sup>18], and connections to non-convex optimization [CDGS20, ZJS20]. Notably, the robust estimators developed in some of these works [DKK<sup>+</sup>16, LRV16, DKK<sup>+</sup>17] are scalable in practice and yield a number of applications in exploratory data analysis [DKK<sup>+</sup>17] and adversarial machine learning [TLM18, DKK<sup>+</sup>18]. The reader is referred to [DK19] for a recent survey.

Our partial clustering algorithm makes essential use of the Sum-of-squares based *proofs to algorithms* framework (see [FKP19] for an exposition). This framework, beginning with [BKS15], uses the *Sum-of-squares method* to design algorithms for statistical estimation problems, and has led to some of the most general outlier-robust learning algorithms. This includes computationally efficient outlier-robust estimators of the mean, covariance, and low-degree moments of structured distributions, with applications to ICA [KS17b], linear regression [KKM18, BP20, ZJS20], clustering spherical mixtures [HL18, KS17a], and clustering non-spherical mixtures [BK20b, DHKK20]. The sum-of-squares method also gives a generally applicable scheme to handle the *list-decodable learning* setting [BBV08, CSV17], where a majority of the input points are corrupted, yielding efficient list-decodable learners for mean estimation [KS17a], regression [KKK19, RY20a], and subspace clustering/recovery [RY20b, BK20a].

Our work also has connections to the usage of tensor decomposition algorithms for learning statistical models. [HK13] used fourth-order tensor decomposition to obtain a polynomial-time algorithm for mixtures of spherical Gaussians with linearly independent means (with condition number guarantees). This result was extended via higher-order tensor decomposition for non-spherical Gaussian mixtures in a smoothed analysis model [GHK15]. Fourth order tensor decomposition has earlier been used in [FJK96] for the ICA problem [FJK96], and extended to general ICA with higher-order tensor decomposition by [GVX14]. Such results rely on additional and non-trivial assumptions on the parameters of the mixture components in order to succeed, and are incomparable to our tensor-decomposition result that does not make any assumptions on the parameters of the mixture components. Indeed, this is the key innovation in our tensor decomposition algorithm that relaxes the guarantees on the output (we output a small list of candidate parameters) under a priori bounds on distance between components that is ensured by our partial clustering subroutine. This relaxation of tensor decomposition, and the new procedure that accomplishes it, is one of the main contributions of our paper.

Finally, we point out that [DKS17] gave an SQ lower bound for learning (fully) clusterable Gaussian mixtures without outliers, which provides evidence that a  $d^{\Omega(k)}$  dependence is necessary in both the sample complexity and runtime of any algorithm that learns GMMs.

## 1.5 Organization

The structure of this paper is as follows: In Section 2, we provide relevant background and technical facts. In Section 3, we describe and analyze our new tensor decomposition algorithm. In Section 4, we use a sum-of-squares based approach to partially cluster a mixture. In Section 5, we give a spectral separation algorithm to identify thin components. In Section 6, we put all these pieces together to prove Theorem 1.4. In Section 7, we present a refinement of our partial clustering procedure that improves the probability of success to a constant independent of the minimum weight of any component in the input mixture. In Section 8, we present an efficient algorithm that replaces an exhaustive search subroutine in the tensor decomposition algorithm and combines it with the improved partial clustering subroutine to get a  $\text{poly}_k(\varepsilon)$ -error guarantee for robust proper learning of Gaussian mixtures and prove Theorem 1.5. Finally, in Section 9, we show that our algorithm in fact achieves the stronger parameter estimation guarantees and prove Theorem 1.6.

## 1.6 Overview of Techniques

### 1.6.1 Proof of Theorem 1.4

In this section, we give a bird’s eye view of our algorithm and the main ideas that go into it. Recall that our goal is to design an efficient algorithm that takes an  $\varepsilon$ -corrupted sample  $Y$  from a mixture of  $k$ -Gaussians  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  and outputs a mixture  $\widehat{\mathcal{M}} = \sum_i \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  such that the total variation distance between  $\mathcal{M}$  and  $\widehat{\mathcal{M}}$  is bounded above by a dimension-independent function of  $\varepsilon$  (bounded above by  $1/(\log(1/\varepsilon))^{k-O(k^2)}$ ). Specifically, we want the running time of our algorithm to be bounded above by a polynomial in the dimension  $d$  and  $1/\varepsilon$ , for any fixed  $k$ .

In the non-robust setting (i.e., for  $\varepsilon = 0$ ), the algorithm of Moitra and Valiant [MV10], extending their work with Kalai [KMV10], solves this problem. However, natural attempts to adapt their method to tolerate outliers run into immediate difficulties. The starting point of [MV10] is to observe that if a mixture of  $k$  Gaussians has every pair of components separated in total variation distance by at least  $\delta$ , then a random univariate projection of the mixture has a pair of components that are  $\delta/\sqrt{d}$ -separated in total variation distance. Their algorithm uses this observation to piece together estimates of the mixture when projected to several carefully chosen directions to get an estimate of the high-dimensional mixture. Notice, however, that such a strategy meets with instant roadblock in the presence of outliers: the fraction of outliers, being a dimension-independent constant, completely overwhelms the total variation distance between components in any one direction making them indistinguishable<sup>2</sup>.

For a reader familiar with the work in algorithmic robust statistics, this may not come as a surprise — to handle outliers, we almost always need to develop a completely new algorithm, even in the outlier-free setting.

In particular, as we next describe, our approach diverges from the method of [MV10] at the very beginning, and instead relies on a careful interleaving of two new algorithmic primitives: (1) a new *partial clustering algorithm* based on the sum-of-squares (SoS) method, and (2) a new *tensor decomposition* method for decomposing a symmetrized sum of tensor powers of  $d \times d$  matrices.

**The Key Determinant: Clusterability.** Our first and key conceptual contribution is to deal with the case of *partially clusterable* mixtures differently from those that are not partially clusterable. We call a mixture partially clusterable if there is a pair of components that have total variation distance larger than  $1 - \Omega_k(1)$  (we will call such components *well-separated* in what follows). We note that even the setting when the mixture is *fully* clusterable (i.e., every pair of components is well-separated), the learning problem captures many hard special cases (such as subspace clustering) and is highly non-trivial. Two recent works [DHKK20, BK20b] gave a polynomial-time algorithm for the fully clusterable case, using the sum-of-squares method. Interestingly, it turns out that the clustering algorithm of [BK20b] (specifically, their Lemma 6.4) can be generalized (see Theorem 4.2) to the partial clustering setting, i.e., the setting where we are guaranteed to have a pair of components that are well-separated (with no guarantees on the remaining components). This gives an algorithm with running time of  $d^{(k/\alpha)^{O(k)}}$  to partition the input sample into components so that each piece of the partition is (effectively) a  $(\text{poly}(\alpha/k) + \varepsilon)$ -corrupted sample from disjoint sub-mixtures. Here,  $\alpha$  is the smallest mixing weight. As we will soon see, such a partial clustering algorithm will be too slow to yield our final guarantees of a fully polynomial time algorithm (when the smallest weight is too small), and we will soon discuss a more efficient variant that will suffice for our purposes.

**Approximate Isotropic Transformation.** By applying our partial-clustering algorithm, we can effectively assume that the input is an  $\varepsilon$ -corrupted sample from a mixture with every pair of

---

<sup>2</sup>Informally speaking, one could hope to show that outliers projected into a random direction cannot be too adversarial, but it is unclear how to use this observation not in the least because the algorithm of [MV10] requires a somewhat carefully tailored choice of projections.

components *at most*  $(1 - \Omega_k(1))$ -far in total variation distance. At this point, we would like to make the mixture isotropic — that is, we would like to assume that the mean of the mixture is  $\approx 0$  and the covariance of the mixture is  $\approx I$ . In the setting with no outliers, this is simply a matter of computing the empirical mean and covariance and applying an appropriate affine transformation to the input points. However, in the setting with outliers, even this task is somewhat non-trivial. A natural idea is to use the algorithm for robust covariance estimation with bounded error in *spectral norm* from [KS17b] that works for all *certifiably subgaussian* distributions (the same work also establishes that arbitrary mixtures of Gaussians are certifiably subgaussian). However, it turns out that our algorithm needs dimension-independent error guarantee on the estimated covariance in *Frobenius norm* (instead of the weaker spectral error guarantee). Fortunately, the recent work [BK20b] (Theorem 7.1 in their paper and Fact 2.38) gives precisely such an algorithm for robust covariance estimation that relies on the stronger property of *certifiable hyper-contractivity* (we verify that this property holds for mixtures of Gaussians in Lemma 2.34 of Section 2).

**Mixtures with Pairwise Close Components.** After the first two steps, we can effectively assume that we are working with an  $\varepsilon$ -corrupted sample from a mixture that is approximately isotropic and every pair of components is not too far in total variation distance. Why is this latter guarantee useful? As established in the recent works [DHKK20, BK20b], such a bound translates into a guarantees (with respect to natural norms) on the parameters of the component Gaussians. In particular, using this translation in our setting implies that after partial clustering plus an approximate isotropic transformation, we have that: 1)  $\|\Sigma_i - I\|_F \leq \text{poly}(\alpha, k)$ , 2)  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$ , and 3)  $\Sigma_i \geq \frac{1}{\text{poly}(\alpha, k)}I$  for every  $i$  (recall that  $\alpha$  is the minimum weight in the mixture). In this case, it turns out that in order to learn the unknown mixture with error guarantees in total variation distance, it suffices to obtain  $\text{poly}_k(\varepsilon)$ -error estimates of the  $\mu_i, \Sigma_i$ 's in Frobenius norm.

**Symmetrized Tensors and the Work of Kane [Kan20].** The key first step in addressing such a mixture was taken in the very recent work of Kane [Kan20], who gave a polynomial-time algorithm to robustly learn an *equiweighted* mixture of two Gaussians. For this special case, after isotropic transformation, one can effectively assume that the two means are  $\pm\mu$  and the two covariances are  $I \pm \Sigma$ . Kane's idea is to look at a certain tensor ("Hermite tensor") that can be built using the 4-th and 6-th raw moments of the mixture. Since we must use outlier-robust algorithms to estimate these tensors, we can obtain estimates that are accurate only up to constant error in Frobenius norm of the tensor. Kane's key observation is that *for the special case of  $k = 2$  components*, one can build two different Hermite tensors, one of which is rank-one with component  $\approx \mu$  (and thus one can immediately "read off"  $\mu$ ); the other only has a tensor power of  $\Sigma$ . This second tensor is of the form  $\hat{T}_4 = \text{Sym}((\Sigma - I) \otimes (\Sigma - I)) + E$ , where  $\|E\|_F = O_k(\sqrt{\varepsilon})$  and  $\text{Sym}$  refers to symmetrizing over all possible permutations of the "4 modes of the tensor". Unlike the case of the mean, one cannot simply "read-off"<sup>3</sup>  $\Sigma$  from  $T_4$ , but Kane gives a simple method to accomplish this. As noted

<sup>3</sup>It is helpful to visualize a single entry of this tensor for, say, the case when  $i, j, k, \ell$  are all distinct:  $\hat{T}_4(i, j, k, \ell) = \frac{1}{3}(\Sigma(i, j)\Sigma(k, \ell) + \Sigma(i, k)\Sigma(j, \ell) + \Sigma(i, \ell)\Sigma(j, k)) + \text{error}$ . Notice that obtaining entries of  $\Sigma$  from  $T_4$  is formally a task of solving noisy quadratic equations.



in [Kan20], it is not clear how to extend this to non-equiweighted mixtures of  $k = 2$  Gaussians, and going to even  $k = 3$  components requires substantially new ideas.

**List-decodable Tensor Decomposition.** Our method for  $k > 2$  works by abstracting and generalizing key aspects of the [Kan20] somewhat ad-hoc approach for the case of  $k = 2$  and combining it with new ideas. This is necessary because of several issues, as we soon discuss: 1) as stated, Kane’s approach does not work as is even for  $k = 2$  when the mixture is not equiweighted, 2) it is not known whether one can build tensors that “separate” out the means and the covariances, as [Kan20] managed to do for  $k = 2$ , 3) the relevant tensors will not be (symmetrizations of) rank-1 tensors, up to noise. It is worth noting, in fact, that such a gap is information-theoretically inherent, as shown in [MV10]; learning arbitrary mixtures of  $k$  Gaussians *provably* requires at least  $2k$  moments of the mixture. Another way of seeing this is by considering the parallel pancakes construction of [DKS17], where the authors produce an example of a mixture of  $k$  Gaussians whose first  $k$  moments match the standard Gaussian exactly despite not being close in total variational distance. It should be noted that, in this example, the component Gaussians are all equal to the standard Gaussian except in one hidden direction. Thus, we cannot hope to identify the components exactly with just  $O(1)$  many moments. (However, for this example at least, we might hope to identify these components up to whatever is going on in this one (unknown) hidden direction. In fact, more complicated constructions can be made to have several such hidden directions.) The somewhat surprising fact (that we establish in more detail below) is that by looking at only the first four moments of our mixture, we can learn all of the components up to some errors taking place along a bounded number of hidden directions. In particular, we can learn the covariance matrices of the component Gaussians up to some *low rank* error terms. We elaborate on this new idea below.

For the sake of the intuition, it is helpful to focus on the simpler case where all the means are zero. In this case, the estimated 4th Hermite tensor (built from estimated raw moments of degree at most 4 of the mixture) has the following form:

$$\hat{T}_4 = \sum_{i=1}^k w_i \text{Sym}((\Sigma_i - I) \otimes (\Sigma_i - I) + E),$$

where  $E$  is a 4-tensor with  $\|E\|_F = O_k(\sqrt{\varepsilon})$ . Given the form of this tensor, it is natural to think of applying tensor decomposition algorithms, by thinking of  $\Sigma_i - I$  as a  $d^2$ -dimensional vector. However, we run into the issue of uniqueness of tensor decomposition, since we are dealing with 2nd order tensors (once we view  $\Sigma_i - I$  as a  $d^2$ -dimensional vector). One might imagine computing higher-order tensors of similar forms to overcome the uniqueness issues, but this runs into two major complications: first, the symmetrization operation introduces spurious terms that do not have the sum of tensor-power structure required for such an algorithm to succeed. Indeed, this is far from just being an annoying technicality — the recent work of Garg, Kayal, and Saha [GKS20] addresses precisely such a tensor decomposition problem via algebraic techniques.

Unfortunately, as is explicitly pointed out as one of the main open question in their work (see page 14), because of their reliance on algebraic techniques, their algorithm is highly brittle and in particular, may not even be able to handle the benign noise that comes from estimating tensors

from independent (uncorrupted) samples. (Of course, our setting has to deal with the malicious noise introduced due to the adversarial outliers.) Second, even if one were to get hold of the tensor without the symmetrization operation, the only applicable tensor decomposition algorithm (recall that we do not make *any* genericity assumptions on the components that are typically required by tensor decomposition algorithms) is the result of Barak, Kelner, and Steurer [BKS15]. However, the [BKS15] result, while being efficient in its dependence on the number of components, has exponential dependence on the target error, which is prohibitively expensive for our application.

Our idea is to give up on the goal of recovering the unique decomposition of the tensor  $\hat{T}_4$  above, and start by applying an operation that is a common trick in most tensor decomposition algorithms. In our context, this trick amounts to taking a random matrix (say, with independent standard Gaussian entries)  $P$  and “collapsing” the last two modes of  $\hat{T}_4$  with  $P$  (i.e., computing  $\hat{S}(i, j) = \sum_{k, \ell} \hat{T}_4(i, j, k, \ell)P(k, \ell)$ ) to obtain a matrix  $Q$ . In the usual tensor decomposition procedures, we are interested in proving that one can recover all the information about the components of the tensor from  $Q$ . We will not be able to prove such a statement here. Instead, our key observation is that one can choose a matrix  $P$  of rank  $\text{poly}(k)$  and argue that the resulting  $\hat{S}$  is  $O_k(\varepsilon)$ -close to one of the  $\Sigma_i - I$  up to an error term of  $O(k^2)$  rank. To see this, note that in the symmetrization

$$\hat{T}_4(i, j, k, \ell) = \frac{1}{3}(\Sigma(i, j) \otimes \Sigma(k, \ell) + \Sigma(i, k) \otimes \Sigma(j, \ell) + \Sigma(i, \ell) \otimes \Sigma(j, k)) + \text{error}$$

applying a rank-one matrix  $P$  to the modes  $k, \ell$  will reduce the first term to the matrix  $\Sigma$ , while the latter two terms become rank-one terms! When the tensor is a sum of  $k$  such symmetrized tensors, we will use  $k$  such rank-one matrices, and take a linear combination of them to get a weighted sum of the  $\Sigma$ 's plus a term of rank  $O(k^2)$ . As we show, the linear combination can be chosen such that only one of the component  $\Sigma$ 's survives (up to small Frobenius norm). Moreover, such a low-rank  $P$  can be obtained efficiently by simply choosing  $\text{poly}(k)$  random rank-1 matrices and exhaustively searching over an appropriate cover of the  $O(k^2)$ -dimensional subspace spanned by them.

**Subspace Enumeration to Recover the Low-rank Terms.** We then show that we can use the generated estimates  $\hat{S}$  and the tensors  $\hat{T}_m$  for  $m \leq 4k$  to find a  $\text{poly}_k(\varepsilon)$ -dimensional subspace  $V_*$  such that the low-rank error matrix in the estimated  $\hat{S}$  have their range space essentially contained inside  $V_*$ . Our next step involves a subspace enumeration over  $V_*$  to output a list of a bounded number of parameters such that a mixture defined by some  $k$  of them must be close in total variation distance to the input mixture.

Our final step involves a relatively standard hypothesis testing procedure, using a robust tournament that goes over each of the candidate mixtures in our generated list and finds one that is approximately the closest in total variation distance to a (fresh) set of corrupted samples.

**An Algorithm with Exponential Dependence in  $1/\varepsilon$ .** The above steps suffice to immediately obtain an algorithm whose running time grows exponentially in the reciprocal of the minimum weight of the mixture. This gives a polynomial-time algorithm (see Theorem 4.12 for details) for robustly learning arbitrary equiweighted mixtures of  $k$  Gaussians. When the weights are not all

equal, notice that we can treat any component with weight  $\leq \varepsilon$  as outliers, which effectively means that  $\alpha \geq \varepsilon$ . Thus, our discussion already yields a  $d^{f(k/\varepsilon)}$ -time algorithm in the general setting.

In order to improve this running time to have a fixed polynomial dependence on  $d$  (independent of  $1/\varepsilon$ ), we rely on a new partial clustering result that weakens the separation guarantee of total variation distance. Our final algorithm then involves a recursive interleaving of the partial clustering and tensor decomposition steps with a new *Recursive Spectral Clustering* subroutine. We discuss these steps next.

**Partial Clustering.** The key bottleneck in the running time guarantee of the algorithm described above is the partial clustering step, so it is important to examine the cause for the exponential dependence on the minimum weight in the running time. The algorithm relies on a recently established characterization of well-separated pair of Gaussians  $\mathcal{N}(\mu_1, \Sigma_1)$  and  $\mathcal{N}(\mu_2, \Sigma_2)$  in terms of three geometric distances between their parameters: 1) Mean Separation:  $\exists v$  such that  $\langle \mu_1 - \mu_2, v \rangle^2 \geq \Delta(v^\top(\Sigma_1 + \Sigma_2)v)$ , 2) Relative Frobenius Separation:  $\left\| \Sigma_1^{-1/2}(\Sigma_2 - \Sigma_1)\Sigma_1^{-1/2} \right\|_F \geq \Delta$ , and 3) Spectral Separation:  $\exists v$  such that  $v^\top \Sigma_1 v \geq \Delta v^\top \Sigma_2 v$  or  $v^\top \Sigma_2 v \geq \Delta v^\top \Sigma_1 v$ . The main idea of the algorithm is to give efficient (low-degree) *sum-of-squares* certificates of *simultaneous intersection bounds* that show that any cluster a natural SoS relaxation finds cannot significantly overlap with two well-separated clusters simultaneously. This step requires sum-of-squares certificates for two natural analytic properties: *certifiable hypercontractivity* and *certifiable anti-concentration* (introduced in the recent works [KKK19, RY20a], and also used in [RY20b, BK20a]). The bottleneck that results in the bad running time for us is the degree of the sum-of-squares certificate needed for certifiable anti-concentration (which grows polynomially in  $1/\alpha$ ).

It is not known whether there is a sum-of-squares certificate of much smaller degree for certifiable anti-concentration. To make progress here, we observe that the only usage of this certificate occurs in dealing with spectrally separated pairs of Gaussians in the mixture. Indeed, we give a new partial clustering algorithm that works in fixed polynomial time, whenever there is a pair of Gaussian components separated either via their means or an appropriate variant of the relative Frobenius distance.

**Tensor Decomposition Needs to be Augmented.** While we gain on the running time through our new partial clustering algorithm, the guarantees of the tensor decomposition subroutine we discussed above are no longer enough to guarantee a recovery of parameters that result in a mixture close in total variation distance. Because of the three conditions that we assumed in the working of the tensor decomposition algorithm, we can no longer guarantee the third one that gives a lower bound on the smallest eigenvalue of every covariance (relative to the covariance of the mixture). In particular, we can end up in a situation where, even though we have a list of parameters that contain Frobenius-norm-close estimates of the covariances, the estimates are not enough to provide a total variation distance guarantee. (Consider for example a “skinny” direction where the variance of some component is very small, or even 0. Then we have to learn the parameters more precisely!)



**Spectral Separation of Thin Components.** It turns out that the above is the only way the algorithm can fail at this point — one or more covariance matrices have a very small eigenvalue (if not, the Frobenius norm error would imply TV-distance error). But since we have estimates of the covariances, we can find such a small eigenvector. Now we observe that since the mixture is nearly isotropic (i.e., the overall variance in each direction is  $\sim 1$ ), if some component has very small variance along a direction, then the components must be separable along this direction. We show that it is possible to efficiently cluster the mixture after projecting it to this direction, so that each cluster has strictly fewer components. We then recursively apply the entire algorithm on the clusters obtained, which will each have strictly fewer components.

**Polynomial Complexity.** To avoid  $\varepsilon$ -dependence in the exponent of the dimension or exponential dependence in the minimum mixing weight, we use our new partial clustering algorithm that does not rely on certifiable anti-concentration, by avoiding spectral proximity guarantees, and moving the work of separating along small eigenvalue directions to later in the algorithm. Our tensor decomposition also has a dependence on the minimum mixing weight. To circumvent this exponential dependence on the minimum weight, a natural approach would be to ignore components lighter than some threshold (that depends on the target error) and treat them as corruptions. However, this intuitive approach runs into a difficulty. In order to get nontrivial error guarantees on the tensor decomposition, we need that the minimum mixing weight is significantly *larger* than the fraction of outliers (since the decomposition involves generating a list of candidate hypotheses, one of which must be accurate). To solve this problem, we show that we can set a minimum weight threshold that depends on the number of remaining components of the mixture, and “remove” some (but not all) components, so that the remaining mixture has minimum weights above this threshold; and the threshold is also sufficiently larger than the total weight of small components treated as corruptions. This makes the overall computational complexity dependence on  $d$  truly polynomial for any fixed  $k$ , and avoids any dependence on the true minimum mixing weight.

### 1.6.2 Proofs of Theorems 1.5 and 1.6

**Robust Proper Learning with  $\text{poly}_k(\varepsilon)$  Error.** The algorithm of Theorem 1.4 achieves error polynomial in  $\varepsilon$ , alas in time exponential in  $\text{poly}_k(1/\varepsilon)$ . This exponential dependence on  $1/\varepsilon$  comes from two sources: (1) the exhaustive subspace enumeration within the space  $V_*$  of significant eigenvalues, and (2) the error probability in the partial clustering algorithm, which then necessitates super-polynomial enumeration. We use additional techniques to reduce both dependencies (and hence the overall running time) to  $\text{poly}_k(1/\varepsilon)$ .

The key bottleneck in the partial clustering step is the rounding algorithm that obtains a list of candidate clusters each of size roughly  $\alpha n$  where  $\alpha$  is the weight of the smallest component cluster (note that this can be arbitrarily small compared to  $1/k$ ). The rounding guarantees that any such candidate cluster cannot simultaneously contain a significant fraction of points from two components whose covariances are separated in Frobenius distance. However, the candidates could be indiscriminate in collecting an arbitrary subset of points from “nearby” clusters. Our rounding

algorithm in Section 4 simply guesses a partition of them such that every pair a  $\exp(-1/\alpha)$  success probability in guessing a 2 partition of points such that any cluster (up to a tiny fraction of errors) appears on one side only.

Our new clustering algorithm improves on this random guessing step in the rounding by observing that points that come from nearby clusters must have covariances that are close in Frobenius norm. Thus our rounding applies a robust covariance estimation algorithm (with error guarantees in Frobenius norm) to each candidate cluster of size roughly  $\alpha n$  and then collects candidate clusters whenever their estimated covariances are close. We show that this procedure coalesces the  $O(1/\alpha)$  clusters into a collection of at most  $k$ . The random guessing step now succeeds with  $\exp(-k)$  probability resolving the bottleneck in the discussion above.

The other bottleneck is in the subspace enumeration over  $V_*$  outputs a list of size with exponential dependence on the dimension of  $V_*$ . To reduce the list size of tensor decomposition, first we observe that we can use an elementary filtering technique to denoise the data, since the number of samples needed is small. Specifically, if the sample complexity of the algorithm is polynomial in the dimension and error parameter, we can set both the dimension of  $V_*$  and the error parameter to be polynomials in  $1/\varepsilon$ , so that the number of samples needed is  $O(1/\varepsilon)$ . Then, with probability  $(1 - \varepsilon)^{O(1/\varepsilon)} = \Omega(1)$ , a sample of size  $O(1/\varepsilon)$  drawn from the total variation contamination model has no noise. For such a clean sample, we can apply a non-robust algorithm to the subspace  $V_*$ . If the time complexity of the algorithm is polynomial in the dimension and error parameter, then the running time to apply it on  $V_*$  will be  $\text{poly}_k(\varepsilon)$ , as desired. The next step is to prove that such an algorithm exists. We will use the algorithm in Theorem 8 of [MV10]. A small technical issue is that the latter assumes that any pair of components in the mixture has TV distance at least  $\delta$ , where  $\delta$  is the error parameter. We show that with an appropriately chosen parameter  $\delta' = \text{poly}(\delta)$ , any pair of components with TV distance less than  $\delta'$  are close enough so that the algorithm cannot distinguish the pair from a single Gaussian, since the algorithm only requires a polynomial number of samples. If we merge all such pairs, any pair in the mixture is separated by  $\delta'$ , and then we can apply Theorem 8 in [MV10]. For each estimate  $\hat{S}$  of the main algorithm, we can recover the low-rank error of  $\hat{S}$  by learning the mixture in the subspace  $V_*$ . Combining the estimates  $\hat{S}$  and the estimates in the subspace  $V_*$ , we get a list of parameters of size  $\text{poly}_k(\varepsilon)$ .

**Parameter Recovery.** We show that for any two Gaussian mixtures, if they are close in TV distance, their components are also close in TV distance, which implies that we can recover the components or the parameters of the mixture. This result generalizes Theorem 8.1 in [LM20], which has three additional assumptions: (i) each component has variance at least  $\text{poly}(\varepsilon/d)$  and at most  $\text{poly}(d/\varepsilon)$ , (ii) each pair of components has TV distance at least  $\text{poly}_k(\varepsilon)$ , and (iii) the minimal weights of both mixtures are at least  $\text{poly}_k(\varepsilon)$ . [LM20] also proved the conclusion under the assumption that any parameters (means and covariances) are identical or separated. We reduce the general case to this simplification. The first step is to deal with the components with small weights. We use a threshold  $\varepsilon' = \text{poly}_k(\varepsilon)$  of weights such that if we treat components with weights smaller than the threshold as noise, other components have weights at least  $\text{poly}_k(\varepsilon')$ . The second step is a partial clustering on the union of the components of the two mixtures, after which the components within each cluster

are pairwise not too close, i.e., have TV distance bounded by  $1 - \text{poly}(\varepsilon')$ . Then we can modify the parameters in each cluster slightly, so that the resulting parameters for different components are either identical or have a minimum separation. Thus, we reduce the general case of two arbitrary mixtures to this special case. Overall, this is a purely information-theoretic statement — for each significant weight component of one mixture, there will be a component in the second mixture with very close mean and covariance. We note that this is not necessarily a 1-1 mapping between components, which is impossible in general without further assumptions.

## 2 Preliminaries

**Basic Notation.** For a vector  $v$ , we use  $\|v\|_2$  to denote its Euclidean norm. For an  $n \times m$  matrix  $M$ , we use  $\|M\|_{\text{op}} = \max_{\|x\|_2=1} \|Mx\|_2$  to denote the operator norm of  $M$  and  $\|M\|_F = \sqrt{\sum_{i,j} M_{i,j}^2}$  to denote the Frobenius norm of  $M$ . We sometimes use the notation  $M(i, j)$  to index the corresponding entries in  $M$ . For an  $n \times n$  symmetric matrix  $M$ , we use  $\geq$  to denote the PSD/Loewner ordering over eigenvalues of  $M$  and  $\text{tr}(M) = \sum_{i \in [n]} M_{i,i}$  to denote the trace of  $M$ . We use  $U\Lambda U^\top$  to denote the eigenvalue decomposition, where  $U$  is an  $n \times n$  matrix with orthonormal columns and  $\Lambda$  is the  $n \times n$  diagonal matrix of the eigenvalues. We use  $M^\dagger = U\Lambda^\dagger U^\top$  to denote the Moore-Penrose pseudoinverse, where  $\Lambda^\dagger$  inverts the non-zero eigenvalues of  $M$ . If  $M \geq 0$ , we use  $M^{+/2} = U\Lambda^{+/2}U^\top$  to denote taking the square-root of the inverted non-zero eigenvalues.

For  $d \times d$  matrices  $A, B$ , the Kronecker product of  $A, B$ , denoted by  $A \otimes B$ , is indexed by  $(i, j), (k, \ell) \in [d] \times [d]$  and has entries  $(A \otimes B)((i, j), (k, \ell)) = A(i, k)B(j, \ell)$ . We will equip every tensor  $T$  with the norm  $\|\cdot\|_F$  that simply corresponds to the  $\ell_2$ -norm of any flattening of  $T$  to a vector. The notation  $T(\cdot, \cdot, x, y)$  is used to denote collapsing two modes of the tensor by plugging in  $x, y$ . For a positive integer  $\ell$  and vector  $v$ , we also use  $v^{\otimes \ell} = \underbrace{v \otimes v \dots \otimes v}_{\ell \text{ times}}$ .

We use the notation  $\mathcal{M} = \sum_{i \in [k]} w_i \mathcal{N}(\mu_i, \Sigma_i)$  to represent a  $k$ -mixture of Gaussians. The total variation distance between two probability distributions on  $\mathbb{R}^d$  with densities  $p, q$  is defined as  $d_{\text{TV}}(p, q) = \frac{1}{2} \int_{\mathbb{R}^d} |p(x) - q(x)| dx$ . We also use  $\mathbb{E}[\cdot]$ ,  $\mathbf{Var}[\cdot]$  and  $\text{Cov}(\cdot)$  to denote the expectation, variance and covariance of a random variable.

For a finite dataset  $X$ , we will use  $Z \in_u X$  to denote that  $Z$  is the uniform distribution on  $X$ . We will sometimes use the term mean (resp. covariance) of  $X$  to refer to  $\mathbb{E}_{Z \in_u X}[Z]$  (resp.  $\text{Cov}_{Z \in_u X}(Z)$ ).

### 2.1 Gaussian Background

The first few facts in this subsection can be found in Kane [Kan20].

**Fact 2.1.** *The total variation distance between two Gaussians  $\mathcal{N}(\mu_1, \Sigma_1)$  and  $\mathcal{N}(\mu_2, \Sigma_2)$  can be bounded above as follows:*

$$d_{\text{TV}}(\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)) = \mathcal{O}\left(\left(\mu_1 - \mu_2\right)^\top \Sigma_1^\dagger (\mu_1 - \mu_2) + \left\| \Sigma_1^{+/2} (\Sigma_2 - \Sigma_1) \Sigma_1^{+/2} \right\|_F\right).$$

**Fact 2.2** (Theorem 2.4 in [Kan20]). Let  $\mathcal{D}$  be a distribution on  $\mathbb{R}^{d \times d}$ , where  $\mathcal{D}$  is supported on the subset of  $\mathbb{R}^{d \times d}$  corresponding to the set of symmetric PSD matrices. Suppose that  $\mathbb{E}[\mathcal{D}] = \Sigma$  and that for any symmetric matrix  $A$  we have that  $\mathbf{Var}[\text{tr}(AX)] = O(\sigma^2 \|\Sigma^{1/2} A \Sigma^{1/2}\|_F^2)$ . Then, for  $\varepsilon \ll \sigma^{-2}$ , there exists a polynomial-time algorithm that given sample access to an  $\varepsilon$ -corrupted set of samples from  $\mathcal{D}$  returns a matrix  $\hat{\Sigma}$  such that with high probability  $\|\Sigma^{-1/2}(\Sigma - \hat{\Sigma})\Sigma^{-1/2}\|_F = O(\sigma\sqrt{\varepsilon})$ .

**Fact 2.3** (Proposition 2.5 in [Kan20]). Let  $G \sim \mathcal{N}(\mu, \Sigma)$  be a Gaussian in  $\mathbb{R}^d$ . Then, we have that

$$\mathbb{E}[G^{\otimes m}](i_1, \dots, i_m) = \sum_{\substack{\text{Partitions } P \text{ of } [m] \\ \text{into sets of size 1 and 2}}} \bigotimes_{\{a,b\} \in P} \Sigma(i_a, i_b) \bigotimes_{\{c\} \in P} \mu(i_c).$$

We will work with the coefficient tensors of  $d$ -dimensional Hermite polynomials:

**Definition 2.4** (Hermite Tensors). Define the degree- $m$  Hermite polynomial tensor as

$$h_m(x) := \sum_{\substack{\text{Partitions } P \text{ of } [m] \\ \text{into sets of size 1 and 2}}} \bigotimes_{\{a,b\} \in P} -I(i_a, i_b) \bigotimes_{\{c\} \in P} x(i_c).$$

We will use the following fact that relates Hermite moments to the raw moments of any distribution.

**Fact 2.5** (Hermite vs Raw Moments, see, e.g., [Her20]). For any real-valued random variable  $u$ , and  $m \in \mathbb{N}$ ,  $\max_{i \leq m} |\mathbb{E} u^i - \mathbb{E}_{z \sim \mathcal{N}(0,1)} z^i| \leq 2^{O(m)} \max_{i \leq m} |\mathbb{E} h_m(u)|$ . Similarly,  $\max_{i \leq m} |\mathbb{E} h_m(u)| \leq 2^{O(m)} \max_{i \leq m} |\mathbb{E} u^i - \mathbb{E}_{z \sim \mathcal{N}(0,1)} z^i|$ .

**Fact 2.6** (Lemma 2.7 in [Kan20]). If  $G \sim \mathcal{N}(\mu, I + \Sigma)$ , then we have that

$$\mathbb{E}[h_m(G)] = \sum_{\substack{\text{Partitions } P \text{ of } [m] \\ \text{into sets of size 1 and 2}}} \bigotimes_{\{a,b\} \in P} \Sigma(i_a, i_b) \bigotimes_{\{c\} \in P} \mu(i_c).$$

**Fact 2.7** (Lemma 2.8 in [Kan20]). If  $G \sim \mathcal{N}(\mu, I + \Sigma)$ , then  $\mathbb{E}[h_m(G) \otimes h_m(G)]$  is equal to

$$\sum_{\substack{\text{Partitions } P \text{ of } [2m] \\ \text{into sets of size 1 and 2 } a,b \text{ in same half of } [2m]}} \bigotimes_{\{a,b\} \in P} \Sigma(i_a, i_b) \bigotimes_{\substack{\{a,b\} \in P \\ a,b \text{ in different halves of } [2m]}} (I + \Sigma)(i_a, i_b) \bigotimes_{\{c\} \in P} \mu(i_c).$$

**Lemma 2.8** (Slight Strengthening of Lemma 5.2 in [Kan20]). For  $G \sim \mathcal{N}(\mu, \Sigma)$ , the covariance matrix of  $h_m(G)$  satisfies:

$$\|\text{Cov}(h_m(G))\|_{op} \leq \|\mathbb{E}[h_m(G) \otimes h_m(G)]\|_{op} = O\left(m(1 + \|\Sigma\|_F + \|\mu\|_2)\right)^{2m}.$$

This follows from the proof of Lemma 5.2 in [Kan20] by noting that the number of terms in the sum is at most  $2^m$  times the number of partitions of  $[2m]$  into sets of size 1 and 2, which is at most  $O(m)^{2m}$ .

Next, we use upper and lower bounds on low-degree polynomials of Gaussian random variables. We defer the proof of the subsequent Lemma to Appendix A.

**Lemma 2.9** (Concentration of low-degree polynomials). *Let  $T$  be a  $d$ -dimensional, degree-4 tensor such that  $\|T\|_F \leq \Delta$  for some  $\Delta > 0$  and let  $x, y \sim \mathcal{N}(0, I)$ . Then, with probability at least  $1 - 1/\text{poly}(d)$ , the following holds:*

$$\|T(\cdot, \cdot, x, y)\|_F^2 \leq O(\log(d)\Delta^2) .$$

Note that for any matrix  $M$ ,  $\langle M, x \otimes y \rangle$ , where  $x, y \sim \mathcal{N}(0, I)$ , is a degree-2 polynomial in Gaussian random variables. As a result, we have the following anti-concentration inequality.

**Lemma 2.10** (Anti-concentration of bi-linear forms, [CW01]). *Let  $M$  be a  $d \times d$  matrix and let  $x, y \sim \mathcal{N}(0, I)$ . Then, for any  $\zeta \in (0, 1)$ , the following holds:*

$$\mathbb{P}\left[\langle M, x \otimes y \rangle^2 \leq \zeta \mathbb{E}\left[\langle M, x \otimes y \rangle^2\right]\right] \leq O(\sqrt{\zeta}) .$$

## 2.2 Sum-of-Squares Proofs and Pseudo-distributions

We refer the reader to the monograph [FKP19] and the lecture notes [BS16] for a detailed exposition of the sum-of-squares method and its usage in average-case algorithm design.

Let  $x = (x_1, x_2, \dots, x_n)$  be a tuple of  $n$  indeterminates and let  $\mathbb{R}[x]$  be the set of polynomials with real coefficients and indeterminates  $x_1, \dots, x_n$ . We say that a polynomial  $p \in \mathbb{R}[x]$  is a *sum-of-squares (sos)* if there exist polynomials  $q_1, \dots, q_r$  such that  $p = q_1^2 + \dots + q_r^2$ .

### 2.2.1 Pseudo-distributions

Pseudo-distributions are generalizations of probability distributions. We can represent a discrete (i.e., finitely supported) probability distribution over  $\mathbb{R}^n$  by its probability mass function  $D: \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $D \geq 0$  and  $\sum_{x \in \text{supp}(D)} D(x) = 1$ . Similarly, we can describe a pseudo-distribution by its mass function by relaxing the constraint  $D \geq 0$  to passing certain low-degree non-negativity tests.

Concretely, a *level- $\ell$  pseudo-distribution* is a finitely-supported function  $D: \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $\sum_x D(x) = 1$  and  $\sum_x D(x)f(x)^2 \geq 0$  for every polynomial  $f$  of degree at most  $\ell/2$ . (Here, the summations are over the support of  $D$ .) A straightforward polynomial-interpolation argument shows that every level- $\infty$ -pseudo distribution satisfies  $D \geq 0$  and is thus an actual probability distribution. We define the *pseudo-expectation* of a function  $f$  on  $\mathbb{R}^n$  with respect to a pseudo-distribution  $D$ , denoted  $\tilde{\mathbb{E}}_{D(x)}f(x)$ , as

$$\tilde{\mathbb{E}}_{D(x)}f(x) = \sum_x D(x)f(x) . \tag{2.1}$$

The degree- $\ell$  moment tensor of a pseudo-distribution  $D$  is the tensor  $\tilde{\mathbb{E}}_{D(x)}(1, x_1, x_2, \dots, x_n)^{\otimes \ell}$ . In particular, the moment tensor has an entry corresponding to the pseudo-expectation of all monomials of degree at most  $\ell$  in  $x$ . The set of all degree- $\ell$  moment tensors of probability distribution is a convex set. Similarly, the set of all degree- $\ell$  moment tensors of degree- $d$  pseudo-distributions is also convex. Unlike moments of distributions, there is an efficient separation oracle for moment tensors of pseudo-distributions.

**Fact 2.11** ([Sho87, Nes00, Las01, Par13]). For any  $n, \ell \in \mathbb{N}$ , the following set has an  $n^{O(\ell)}$ -time weak separation oracle (in the sense of [GLS81]):

$$\{\tilde{\mathbb{E}}_{D(x)}(1, x_1, x_2, \dots, x_n)^{\otimes d} \mid \text{degree-}d \text{ pseudo-distribution } D \text{ over } \mathbb{R}^n\} . \quad (2.2)$$

This fact, together with the equivalence of weak separation and optimization [GLS81], allows us to efficiently optimize over pseudo-distributions (approximately) — this algorithm is referred to as the sum-of-squares algorithm. The *level- $\ell$  sum-of-squares algorithm* optimizes over the space of all level- $\ell$  pseudo-distributions that satisfy a given set of polynomial constraints (defined below).

**Definition 2.12** (Constrained pseudo-distributions). Let  $D$  be a level- $\ell$  pseudo-distribution over  $\mathbb{R}^n$ . Let  $\mathcal{A} = \{f_1 \geq 0, f_2 \geq 0, \dots, f_m \geq 0\}$  be a system of  $m$  polynomial inequality constraints. We say that  $D$  satisfies the system of constraints  $\mathcal{A}$  at degree  $r$ , denoted  $D \stackrel{|}{\vDash}_r \mathcal{A}$ , if for every  $S \subseteq [m]$  and every sum-of-squares polynomial  $h$  with  $\deg h + \sum_{i \in S} \max\{\deg f_i, r\}$ , we have that  $\tilde{\mathbb{E}}_D h \cdot \prod_{i \in S} f_i \geq 0$ .

We write  $D \stackrel{|}{\vDash} \mathcal{A}$  (without specifying the degree) if  $D \stackrel{|}{\vDash}_0 \mathcal{A}$  holds. Furthermore, we say that  $D \stackrel{|}{\vDash}_r \mathcal{A}$  holds *approximately* if the above inequalities are satisfied up to an error of  $2^{-n^\ell} \cdot \|h\| \cdot \prod_{i \in S} \|f_i\|$ , where  $\|\cdot\|$  denotes the Euclidean norm<sup>4</sup> of the coefficients of a polynomial in the monomial basis.

We remark that if  $D$  is an actual (discrete) probability distribution, then we have that  $D \stackrel{|}{\vDash} \mathcal{A}$  if and only if  $D$  is supported on solutions to the constraints  $\mathcal{A}$ . We say that a system  $\mathcal{A}$  of polynomial constraints is *explicitly bounded* if it contains a constraint of the form  $\{\|x\|^2 \leq M\}$ . The following fact is a consequence of Fact 2.11 and [GLS81]:

**Fact 2.13** (Efficient Optimization over Pseudo-distributions). There exists an  $(n + m)^{O(\ell)}$ -time algorithm that, given any explicitly bounded and satisfiable system<sup>5</sup>  $\mathcal{A}$  of  $m$  polynomial constraints in  $n$  variables, outputs a level- $\ell$  pseudo-distribution that satisfies  $\mathcal{A}$  approximately.

**Basic Facts about Pseudo-Distributions.** We will use the following Cauchy-Schwarz inequality for pseudo-distributions:

**Fact 2.14** (Cauchy-Schwarz for Pseudo-distributions). Let  $f, g$  be polynomials of degree at most  $d$  in indeterminate  $x \in \mathbb{R}^d$ . Then, for any degree- $d$  pseudo-distribution  $\tilde{\zeta}$ , we have that  $\tilde{\mathbb{E}}_{\tilde{\zeta}}[fg] \leq \sqrt{\tilde{\mathbb{E}}_{\tilde{\zeta}}[f^2]} \sqrt{\tilde{\mathbb{E}}_{\tilde{\zeta}}[g^2]}$ .

**Fact 2.15** (Hölder's Inequality for Pseudo-Distributions). Let  $f, g$  be polynomials of degree at most  $d$  in indeterminate  $x \in \mathbb{R}^d$ . Fix  $t \in \mathbb{N}$ . Then, for any degree- $dt$  pseudo-distribution  $\tilde{\zeta}$ , we have that  $\tilde{\mathbb{E}}_{\tilde{\zeta}}[f^{t-1}g] \leq (\tilde{\mathbb{E}}_{\tilde{\zeta}}[f^t])^{\frac{t-1}{t}} (\tilde{\mathbb{E}}_{\tilde{\zeta}}[g^t])^{1/t}$ .

**Corollary 2.16** (Comparison of Norms). Let  $\tilde{\zeta}$  be a degree- $t^2$  pseudo-distribution over a scalar indeterminate  $x$ . Then, we have that  $\tilde{\mathbb{E}}[x^t]^{1/t} \geq \tilde{\mathbb{E}}[x^{t'}]^{1/t'}$  for every  $t' \leq t$ .

<sup>4</sup>The choice of norm is not important here because the factor  $2^{-n^\ell}$  swamps the effects of choosing another norm.

<sup>5</sup>Here, we assume that the bit complexity of the constraints in  $\mathcal{A}$  is  $(n + m)^{O(1)}$ .

### 2.2.2 Sum-of-squares proofs

Let  $f_1, f_2, \dots, f_r$  and  $g$  be multivariate polynomials in  $x$ . A *sum-of-squares proof* that the constraints  $\{f_1 \geq 0, \dots, f_m \geq 0\}$  imply the constraint  $\{g \geq 0\}$  consists of polynomials  $(p_S)_{S \subseteq [m]}$  such that

$$g = \sum_{S \subseteq [m]} p_S \cdot \prod_{i \in S} f_i. \quad (2.3)$$

We say that this proof has *degree*  $\ell$  if for every set  $S \subseteq [m]$  the polynomial  $p_S \prod_{i \in S} f_i$  has degree at most  $\ell$ . If there is a degree  $\ell$  SoS proof that  $\{f_i \geq 0 \mid i \leq r\}$  implies  $\{g \geq 0\}$ , we write:

$$\{f_i \geq 0 \mid i \leq r\} \Big|_{\ell} \{g \geq 0\}. \quad (2.4)$$

For all polynomials  $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$  and for all functions  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m, G: \mathbb{R}^n \rightarrow \mathbb{R}^k, H: \mathbb{R}^p \rightarrow \mathbb{R}^n$  such that each of the coordinates of the outputs are polynomials of the inputs, we have the following inference rules.

The first one derives new inequalities by addition or multiplication:

$$\frac{\mathcal{A} \Big|_{\ell} \{f \geq 0, g \geq 0\}}{\mathcal{A} \Big|_{\ell} \{f + g \geq 0\}}, \frac{\mathcal{A} \Big|_{\ell} \{f \geq 0\}, \mathcal{A} \Big|_{\ell'} \{g \geq 0\}}{\mathcal{A} \Big|_{\ell+\ell'} \{f \cdot g \geq 0\}}. \quad (2.5)$$

The next one derives new inequalities by transitivity:

$$\frac{\mathcal{A} \Big|_{\ell} \mathcal{B}, \mathcal{B} \Big|_{\ell'} C}{\mathcal{A} \Big|_{\ell+\ell'} C}. \quad (2.6)$$

Finally, the last rule derives new inequalities via substitution:

$$\frac{\{F \geq 0\} \Big|_{\ell} \{G \geq 0\}}{\{F(H) \geq 0\} \Big|_{\ell \cdot \deg(H)} \{G(H) \geq 0\}}. \quad (\text{substitution})$$

Low-degree sum-of-squares proofs are sound and complete if we take low-level pseudo-distributions as models. Concretely, sum-of-squares proofs allow us to deduce properties of pseudo-distributions that satisfy some constraints.

**Fact 2.17 (Soundness).** *If  $D \Big|_{\overline{r}} \mathcal{A}$  for a level- $\ell$  pseudo-distribution  $D$  and there exists a sum-of-squares proof  $\mathcal{A} \Big|_{\overline{r'}} \mathcal{B}$ , then  $D \Big|_{\overline{r+r'}} \mathcal{B}$ .*

If the pseudo-distribution  $D$  satisfies  $\mathcal{A}$  only approximately, soundness continues to hold if we require an upper bound on the bit-complexity of the sum-of-squares  $\mathcal{A} \Big|_{\overline{r'}} \mathcal{B}$  (i.e., the number of bits required to write down the proof). In our applications, the bit complexity of all sum-of-squares proofs will be  $n^{O(\ell)}$  (assuming that all numbers in the input have bit complexity  $n^{O(1)}$ ). This bound suffices in order to argue about pseudo-distributions that satisfy polynomial constraints approximately.

The following fact shows that every property of low-level pseudo-distributions can be derived by low-degree sum-of-squares proofs.



**Fact 2.18** (Completeness). Suppose that  $d \geq r' \geq r$  and  $\mathcal{A}$  is a collection of polynomial constraints with degree at most  $r$ , and  $\mathcal{A} \mid \{ \sum_{i=1}^n x_i^2 \leq B \}$  for some finite  $B$ . Let  $\{g \geq 0\}$  be a polynomial constraint. If every degree- $d$  pseudo-distribution that satisfies  $D \stackrel{r}{=} \mathcal{A}$  also satisfies  $D \stackrel{r'}{=} \{g \geq 0\}$ , then for every  $\varepsilon > 0$ , there is a sum-of-squares proof  $\mathcal{A} \mid_d \{g \geq -\varepsilon\}$ .

**Basic Sum-of-Squares Proofs.** We will require the following basic SoS proofs.

**Fact 2.19** (Operator norm Bound). Let  $A$  be a symmetric  $d \times d$  matrix and  $v$  be a vector in  $\mathbb{R}^d$ . Then, we have that

$$\mid_v \{v^\top A v \leq \|A\|_2 \|v\|_2^2\} .$$

**Fact 2.20** (SoS Hölder's Inequality). Let  $f_i, g_i$ , for  $1 \leq i \leq s$ , be scalar-valued indeterminates. Let  $p$  be an even positive integer. Then,

$$\mid_{\frac{f, g}{p^2}} \left\{ \left( \frac{1}{s} \sum_{i=1}^s f_i g_i^{p-1} \right)^p \leq \left( \frac{1}{s} \sum_{i=1}^s f_i^p \right) \left( \frac{1}{s} \sum_{i=1}^s g_i^p \right)^{p-1} \right\} .$$

Observe that using  $p = 2$  above yields the SoS Cauchy-Schwarz inequality.

**Fact 2.21** (SoS Almost Triangle Inequality). Let  $f_1, f_2, \dots, f_r$  be indeterminates. Then, we have that

$$\mid_{\frac{f_1, f_2, \dots, f_r}{2t}} \left\{ \left( \sum_{i \leq r} f_i \right)^{2t} \leq r^{2t-1} \left( \sum_{i=1}^r f_i^{2t} \right) \right\} .$$

**Fact 2.22** (SoS AM-GM Inequality, see Appendix A of [BKS15]). Let  $f_1, f_2, \dots, f_m$  be indeterminates. Then, we have that

$$\mid_{\frac{f_1, f_2, \dots, f_m}{m}} \left\{ \left( \frac{1}{m} \sum_{i=1}^m f_i \right)^m \geq \prod_{i \leq m} f_i \right\} .$$

We defer the proofs of the two subsequent lemmas to Appendix A.

**Lemma 2.23** (Spectral SoS Proofs). Let  $A$  be a  $d \times d$  matrix. Then for  $d$ -dimensional vector-valued indeterminate  $v$ , we have:

$$\mid_v \{v^\top A v \leq \|A\|_2 \|v\|_2^2\} .$$

**Fact 2.24** (Cancellation within SoS, Lemma 9.2 [BK20b]). Let  $a, C$  be scalar-valued indeterminates. Then,

$$\{a \geq 0\} \cup \{a^t \leq C a^{t-1}\} \mid_{\frac{a, C}{2t}} \{a^{2t} \leq C^{2t}\} .$$

**Lemma 2.25** (Frobenius Norms of Products of Matrices). Let  $B$  be a  $d \times d$  matrix valued indeterminate for some  $d \in \mathbb{N}$ . Then, for any  $0 \leq A \leq I$ ,

$$\mid_B \{ \|AB\|_F^2 \leq \|B\|_F^2 \} ,$$

and,

$$\mid_B \{ \|BA\|_F^2 \leq \|B\|_F^2 \} ,$$



### 2.3 Analytic Properties of Gaussian Distributions

The following definitions and results describe the analytic properties of Gaussian distributions that we will use. We also state the guarantees of known robust estimation algorithms for estimating the mean, covariance and moment tensors of Gaussian mixtures here.

**Certifiable Subgaussianity.** We will make essential use of the following definition.

**Definition 2.26** (Certifiable Subgaussianity (Definition 5.1 in [KS17b])). For  $t \in \mathbb{N}$  and an absolute constant  $C > 0$ , a distribution  $\mathcal{D}$  on  $\mathbb{R}^d$  is said to be  $t$ -certifiably  $C$ -subgaussian if for every even  $t' \leq t$ , we have that

$$\frac{v}{t'} \left\{ \mathbb{E}_{\mathcal{D}} \langle x, v \rangle^{t'} \leq (Ct')^{t'/2} \left( \mathbb{E}_{\mathcal{D}} \langle x, v \rangle^2 \right)^{t'/2} \right\}.$$

**Fact 2.27** (Mixtures of Certifiably Subgaussian Distributions, Analogous to Lemma 5.4 in [KS17b]). Let  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_q$  be  $t$ -certifiably  $C$ -subgaussian distributions on  $\mathbb{R}^d$ . Let  $p_1, p_2, \dots, p_q$  be non-negative weights such that  $\sum_i p_i = 1$  and  $p = \min_{i \leq q} p_i$ . Then, the mixture  $\sum_i p_i \mathcal{D}_i$  is  $t$ -certifiably  $C/p$ -subgaussian.

**Certifiable Anti-Concentration.** The first is *certifiable anti-concentration* — an SoS formulation of classical anti-concentration inequalities — that was introduced in [KKK19, RY20a].

In order to formulate certifiable anti-concentration, we start with a univariate even polynomial  $p$  that serves as a uniform approximation to the delta function at 0 in an interval around 0. Such polynomials are constructed in [KKK19, RY20a] (see also [DGJ<sup>+</sup>10]). Let  $q_{\delta, \Sigma}(x, v)$  be a multivariate (in  $v$ ) polynomial defined by  $q_{\delta, \Sigma}(x, v) = (v^\top \Sigma v)^{2s} p_{\delta, \Sigma} \left( \frac{\langle x, v \rangle}{\sqrt{v^\top \Sigma v}} \right)$ . Since  $p_{\delta, \Sigma}$  is an even polynomial,  $q_{\delta, \Sigma}$  is a polynomial in  $v$ .

**Definition 2.28** (Certifiable Anti-Concentration). A mean-0 distribution  $D$  with covariance  $\Sigma$  is  $2s$ -certifiably  $(\delta, C\delta)$ -anti-concentrated if for  $q_{\delta, \Sigma}(x, v)$  defined above, there exists a degree- $2s$  sum-of-squares proof of the following two unconstrained polynomial inequalities in indeterminate  $v$ :

$$\left\{ \langle x, v \rangle^{2s} + \delta^{2s} q_{\delta, \Sigma}(x, v)^2 \geq \delta^{2s} (v^\top \Sigma v)^{2s} \right\}, \left\{ \mathbb{E}_{x \sim D} q_{\delta, \Sigma}(x, v)^2 \leq C\delta (v^\top \Sigma v)^{2s} \right\}.$$

An isotropic subset  $X \subseteq \mathbb{R}^d$  is  $2s$ -certifiably  $(\delta, C\delta)$ -anti-concentrated if the uniform distribution on  $X$  is  $2s$ -certifiably  $(\delta, C\delta)$ -anti-concentrated.

*Remark 2.29.* The function  $s(\delta)$  can be taken to be  $O(\frac{1}{\delta^2})$  for standard Gaussian distribution and the uniform distribution on the unit sphere (see [KKK19] and [BK20a]).

**Certifiable Hypercontractivity.** Next, we define *certifiable hypercontractivity* of degree-2 polynomials that formulates (within SoS) the fact that higher moments of degree-2 polynomials of distributions (such as Gaussians) can be bounded in terms of appropriate powers of their 2nd moment.

**Definition 2.30** (Certifiable Hypercontractivity). An isotropic distribution  $\mathcal{D}$  on  $\mathbb{R}^d$  is said to be  $h$ -certifiably  $C$ -hypercontractive if there is a degree- $h$  sum-of-squares proof of the following unconstrained polynomial inequality in  $d \times d$  matrix-valued indeterminate  $Q$ :

$$\mathbb{E}_{x \sim \mathcal{D}} (x^\top Q x)^h \leq (Ch)^h \left( \mathbb{E}_{x \sim \mathcal{D}} (x^\top Q x)^2 \right)^{h/2}.$$

A set of points  $X \subseteq \mathbb{R}^d$  is said to be  $C$ -certifiably hypercontractive if the uniform distribution on  $X$  is  $h$ -certifiably  $C$ -hypercontractive.

Hypercontractivity is an important notion in high-dimensional probability and analysis on product spaces [O'D14]. Kauters, O'Donnell, Tan and Zhou [KOTZ14] showed certifiable hypercontractivity of Gaussians and more generally product distributions with subgaussian marginals. Certifiable hypercontractivity strictly generalizes the better known *certifiable subgaussianity* property (studied first in [KS17b]) that controls higher moments of linear polynomials.

Observe that the definition above is affine invariant. In particular, we immediately obtain:

**Fact 2.31.** *Given  $t \in \mathbb{N}$ , if a random variable  $x$  on  $\mathbb{R}^d$  has  $t$ -certifiably  $C$ -hypercontractive degree-2 polynomials, then so does  $Ax$  for any  $A \in \mathbb{R}^{d \times d}$ .*

As observed in [KS17b], the Gaussian distribution is  $t$ -certifiably 1-subgaussian and  $t$ -certifiably 1-hypercontractive for every  $t$ . Next, we establish certifiable hypercontractivity for mixtures of Gaussians. We defer the proofs to Appendix A.

**Lemma 2.32** (Shifts Cannot Decrease Variance). *Let  $\mathcal{D}$  be a distribution on  $\mathbb{R}^d$ ,  $Q$  be a  $d \times d$  matrix-valued indeterminate, and  $C$  be a scalar-valued indeterminate. Then, we have that*

$$\left| \frac{Q, C}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x) - \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] \right)^2 \right] \leq \mathbb{E}_{x \sim \mathcal{D}} \left[ (Q(x) - C)^2 \right] \right\} \right|.$$

**Lemma 2.33** (Shifts of Certifiably Hypercontractive Distributions). *Let  $x$  be a mean-0 random variable with distribution  $\mathcal{D}$  on  $\mathbb{R}^d$  with  $t$ -certifiably  $C$ -hypercontractive degree-2 polynomials. Then, for any fixed constant vector  $c \in \mathbb{R}^d$ , the random variable  $x + c$  also has  $t$ -certifiably  $4C$ -hypercontractive degree-2 polynomials.*

**Lemma 2.34** (Mixtures of Certifiably Hypercontractive Distributions). *Let  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$  have  $t$ -certifiably  $C$ -hypercontractive degree-2 polynomials on  $\mathbb{R}^d$ , for some fixed constant  $C$ . Then, any mixture  $\mathcal{D} = \sum_i w_i \mathcal{D}_i$  also has  $t$ -certifiably  $(C/\alpha)$ -hypercontractive degree-2 polynomials for  $\alpha = \min_{i \leq k, w_i > 0} w_i$ .*

**Corollary 2.35** (Certifiable Hypercontractivity of Mixtures of  $k$  Gaussians). *Let  $\mathcal{M}$  be a  $k$ -mixture of Gaussians  $\sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  with weights  $w_i \geq \alpha$  for every  $i \in [k]$ . Then, for all  $t \in \mathbb{N}$ ,  $\mathcal{D}$  has  $t$ -certifiably  $4/\alpha$ -hypercontractive degree-2 polynomials.*

We will use the following robust mean estimation algorithm for bounded covariance distributions [DKK<sup>+</sup>17]:

**Fact 2.36** (Robust Mean Estimation for Bounded Covariance Distributions). *There is a poly( $n$ ) time algorithm that takes input an  $\varepsilon$ -corruption  $Y$  of a collection of  $n$  points  $X \subseteq \mathbb{R}^d$ , and outputs an estimate  $\hat{\mu}$  satisfying  $\|\mathbb{E}_{x \sim_u X} x - \hat{\mu}\|_2 \leq O(\sqrt{\varepsilon}) \|\mathbb{E}_{x \sim_u X} (x - \mathbb{E}_{x \sim_u X} x)(x - \mathbb{E}_{x \sim_u X} x)^\top\|_2$ .*

We will use the following robust covariance estimation algorithm from [KS17b]:

**Fact 2.37** (Robust Covariance Estimation, [KS17b]). *For every  $C > 0$ ,  $\varepsilon > 0$  and even  $k \in \mathbb{N}$  such that  $Ck\varepsilon^{1-2/k} \leq c$  for some small enough absolute constant  $c$ , there exists a polynomial-time algorithm that given an (corrupted) sample  $S$  outputs an estimate of the covariance  $\hat{\Sigma} \in \mathbb{R}^{d \times d}$  with the following guarantee: there exists  $n_0 \geq (C + d)^{O(k)} / \varepsilon$  such that if  $S$  is an  $\varepsilon$ -corrupted sample with size  $|S| \geq n_0$  of a  $k$ -certifiably  $C$ -subgaussian distribution  $D$  over  $\mathbb{R}^d$  with mean  $\mu \in \mathbb{R}^d$  and covariance  $\Sigma \in \mathbb{R}^{d \times d}$ , then with high probability:*

$$(1 - \delta)\Sigma \geq \hat{\Sigma} \geq (1 + \delta)\Sigma$$

for  $\delta \leq O(Ck)\varepsilon^{1-2/k}$ .

We will also require the following robust estimation algorithm with Frobenius distance guarantees proven for certifiably hypercontractive distributions in [BK20b]. Since we obtain estimates to the true covariance in Lowner ordering, we can obtain the subspace spanned by the inliers exactly, project on to this subspace and apply Theorem 7.1 in [BK20b].

**Fact 2.38** (Robust Mean and Covariance Estimation for Certifiably Hypercontractive Distributions, Theorem 7.1 in [BK20b]). *Given  $t \in \mathbb{N}$ , and  $\varepsilon > 0$  sufficiently small so that  $Ct\varepsilon^{1-4/t} \ll 1$ <sup>6</sup>, for some absolute constant  $C > 0$ . Then, there is an algorithm that takes input  $Y$ , an  $\varepsilon$ -corruption of a sample  $X$  of size  $n$  with mean  $\mu_*$ , covariance  $\Sigma_*$ , and  $2t$ -certifiably  $C$ -hypercontractive degree-2 polynomials, runs in time  $n^{O(t)}$ , and outputs an estimate  $\hat{\mu}$  and  $\hat{\Sigma}$  satisfying:*

1.  $\left\| \Sigma_*^{t/2} (\mu_* - \hat{\mu}) \right\|_2 \leq O(Ct)^{1/2} \varepsilon^{1-1/t}$ ,
2.  $(1 - \eta)\Sigma_* \leq \hat{\Sigma} \leq (1 + \eta)\Sigma_*$  for  $\eta \leq O(Ck)\varepsilon^{1-2/t}$ , and,
3.  $\left\| \Sigma_*^{t/2} (\hat{\Sigma} - \Sigma_*) \Sigma_*^{t/2} \right\|_F \leq (C't)O(\varepsilon^{1-1/t})$ ,

where  $C' = \max\{C, B\}$  for the smallest possible  $B > 0$  such that for  $d \times d$ -matrix-valued indeterminate  $Q$ ,  $\left\{ \frac{Q}{2} \left\{ \mathbb{E}_{\mathcal{D}} (x^\top Qx - \mathbb{E}_{\mathcal{D}} x^\top Qx)^2 \leq B \left\| \Sigma_*^{1/2} Q \Sigma_*^{1/2} \right\|_F^2 \right\} \right\}$ <sup>7</sup>.

The last line in the above fact asserts a bound (along with a degree 2 SoS proof) on the variance of degree 2 polynomials in terms of the Frobenius norm of its coefficient matrix. In the next few claims, we verify this property via elementary arguments for the two classes of distributions relevant to

<sup>6</sup>This notation means that we needed  $Ct\varepsilon^{1-2/t}$  to be at most  $c_0$  for some absolute constant  $c_0 > 0$ .

<sup>7</sup>The first two guarantees here hold for the larger class of certifiably subgaussian distributions and were proven in [KS17b] (see Theorem 1.2). Gaussian distribution (with arbitrary mean and covariance) are  $t$ -certifiably 1-subgaussian for all  $t$  and their mixtures (similar to Lemma 2.34 and explicitly proven in Lemma 5.4 of [KS17b]) are  $t$ -certifiably  $O(1/\alpha)$ -subgaussian where  $\alpha$  is the minimum mixing weight.

this paper. We note that whenever a distribution satisfies the bounded variance property (without an SoS proof), it also satisfies the property via a degree 2 SoS proof using Lemma 2.23. Thus, asking for an SoS proof of degree 2 in this context poses no additional restrictions on the distribution. Nevertheless, we provide explicit and direct SoS proofs in the following.

We first note that this property of having *certifiable bounded variance* is closed under linear transformations.

**Lemma 2.39** (Linear Transformations of Certifiably Bounded-Variance Distributions). *For  $d \in \mathbb{N}$ , let  $x$  be a random variable with distribution  $\mathcal{D}$  on  $\mathbb{R}^d$  such that for  $d \times d$  matrix-valued indeterminate  $Q$ ,  $\frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} (x^\top Q x - \mathbb{E}_{\mathcal{D}} x^\top Q x)^2 \leq \left\| \Sigma^{1/2} Q \Sigma^{1/2} \right\|_F^2 \right\}$ . Let  $A$  be an arbitrary  $d \times d$  matrix and let  $x' = Ax$  be the random variable with covariance  $\Sigma' = AA^\top$ . Then, we have that*

$$\frac{Q}{2} \left\{ \mathbb{E}_{x' \sim \mathcal{D}'} (x'^\top Q x' - \mathbb{E}_{\mathcal{D}'} x'^\top Q x')^2 \leq \left\| \Sigma'^{1/2} Q \Sigma'^{1/2} \right\|_F^2 \right\}.$$

**Lemma 2.40** (Variance of Degree-2 Polynomials of Standard Gaussians). *We have that*

$$\frac{Q}{2} \left\{ \mathbb{E}_{N(0,I)} \left( x^\top Q x - \mathbb{E}_{N(0,I)} x^\top Q x \right)^2 \leq 3 \|Q\|_F^2 \right\}.$$

*Remark 2.41.* As is easy to verify, the same proof more generally holds for any distribution that has the same first four moments as the zero-mean Gaussian distribution.

As an immediate corollary of the previous two lemmas, we have:

**Corollary 2.42** (Variance of Degree-2 Polynomials of Zero-Mean, Arbitrary Covariance Gaussians). *For any  $0 \preceq \Sigma$ , we have that*

$$\frac{Q}{2} \left\{ \mathbb{E}_{N(0,\Sigma)} \left( x^\top Q x - \mathbb{E}_{N(0,\Sigma)} x^\top Q x \right)^2 \leq 3 \left\| \Sigma^{1/2} Q \Sigma^{1/2} \right\|_F^2 \right\}.$$

We next prove that the same property holds for mixtures of Gaussians satisfying certain conditions.

**Lemma 2.43** (Variance of Degree-2 Polynomials of Mixtures). *Let  $\mathcal{M} = \sum_i w_i \mathcal{D}_i$  be a  $k$ -mixture of distributions  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$  with means  $\mu_i$  and covariances  $\Sigma_i$ . Let  $\mu = \sum_i w_i \mu_i$  be the mean of  $\mathcal{M}$ . Suppose that each of  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$  have certifiably  $C$ -bounded-variance i.e. for  $Q$ : a symmetric  $d \times d$  matrix-valued indeterminate.*

$$\frac{Q}{2} \left\{ \mathbb{E}_{x' \sim \mathcal{D}_i} (x'^\top Q x' - \mathbb{E}_{\mathcal{D}_i} x'^\top Q x')^2 \leq C \left\| \Sigma_i^{1/2} Q \Sigma_i^{1/2} \right\|_F^2 \right\}.$$

*Further, suppose that for some  $H > 1$ ,  $\|\mu_i - \mu\|_2^2, \|\Sigma_i - I\|_F \leq H$  for every  $1 \leq i \leq k$ . Then, we have that*

$$\frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{M}} \left[ \left( x^\top Q x - \mathbb{E}_{x \sim \mathcal{M}} [x^\top Q x] \right)^2 \right] \leq 100CH^2 \|Q\|_F^2 \right\}.$$

As an immediate corollary of Lemma 2.39 and Lemma 2.43, we obtain:

**Lemma 2.44** (Variance of Degree-2 Polynomials of Mixtures of Gaussians). *Let  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with  $w_i \geq \alpha$ , mean  $\mu = \sum_i w_i \mu_i$  and covariance  $\Sigma = \sum_i w_i ((\mu_i - \mu)(\mu_i - \mu)^\top + \Sigma_i)$ . Suppose that for some  $H > 1$ ,  $\|\Sigma^{+1/2}(\Sigma_i - \Sigma)\Sigma^{+1/2}\|_F \leq H$  for every  $1 \leq i \leq k$ . Let  $Q$  be a symmetric  $d \times d$  matrix-valued indeterminate. Then for  $H' = \max\{H, 1/\alpha\}$ ,*

$$\frac{|Q|}{2} \left\{ \mathbb{E}_{x \sim \mathcal{M}} \left[ \left( x^\top Q x - \mathbb{E}_{x \sim \mathcal{M}} [x^\top Q x] \right)^2 \right] \leq 100H'^2 \|\Sigma^{1/2} Q \Sigma^{1/2}\|_F^2 \right\}.$$

**Analytic Properties are Inherited by Samples.** The following lemma can be proven via similar, standard techniques as in several prior works [KS17b, KKK19, BK20a, BK20b].

**Fact 2.45.** *Let  $D$  be a distribution on  $\mathbb{R}^d$  with mean  $\mu$  and covariance  $\Sigma$ . Let  $t \in \mathbb{N}$ . Let  $X$  be a sample from  $D$  such that,  $\left\| \frac{1}{|X|} \sum_{x \in X} (1, \bar{x})^{\otimes t} - \mathbb{E}_{x \sim D} (1, \bar{x})^{\otimes t} \right\|_F \leq d^{-O(t)}$ . Here,  $\bar{x} = \Sigma^{+1/2}(x - \mu)$ . Then,*

1. *If  $D$  is  $2t$ -certifiably  $C$ -subgaussian, then the uniform distribution on  $X$  is  $t$ -certifiably  $2C$ -subgaussian.*
2. *If  $D$  has  $2t$ -certifiably  $C$ -hypercontractive degree 2 polynomials, then the uniform distribution on  $X$  has  $t$ -certifiably  $2C$ -hypercontractive degree 2 polynomials.*
3. *If  $D$  is  $2t$ -certifiably  $C\delta$ -anti-concentrated, then the uniform distribution on  $X$  is  $t$ -certifiably  $2C\delta$ -anti-concentrated.*
4. *If  $\frac{|Q|}{2} \{ \mathbb{E}_{x \sim \mathcal{D}} (x^\top Q x - \mathbb{E}_{x \sim \mathcal{D}} x^\top Q x)^2 \leq C \|Q\|_F^2 \}$ , then, for the uniform distribution  $\mathcal{D}_X$  on  $X$ ,  $\frac{|Q|}{2} \{ \mathbb{E}_{x \sim \mathcal{D}_X} (x^\top Q x - \mathbb{E}_{x \sim \mathcal{D}_X} x^\top Q x)^2 \leq 2C \|Q\|_F^2 \}$ .*

## 2.4 Deterministic Conditions on the Uncorrupted Samples

In this section, we describe the set of deterministic conditions on the set of uncorrupted samples, under which our algorithms succeed. We will require the following definition.

**Definition 2.46.** Fix  $0 < \varepsilon < 1/2$ . We say that a multiset  $Y$  of points in  $\mathbb{R}^d$  is an  $\varepsilon$ -corrupted version (or an  $\varepsilon$ -corruption) of a multiset  $X$  of points in  $\mathbb{R}^d$  if  $|X \cap Y| \geq \max\{(1 - \varepsilon)|X|, (1 - \varepsilon)|Y|\}$ .

Throughout this paper and unless otherwise specified, we will use  $X$  to denote a multiset of i.i.d. samples from the target  $k$ -mixture  $\mathcal{M} = \sum_{i=1}^k w_i G_i$ , where  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$ . We will use  $X_i$  for the subset of points in  $X$  drawn from  $G_i$ , i.e.,  $X = \cup_{i=1}^k X_i$ .

We will use  $Y$  to denote an  $\varepsilon$ -corrupted version of  $X$ , as per Definition 2.46. In this *strong contamination model*, the adversary can see the clean samples from  $X$  before they decide on the  $\varepsilon$ -corruption  $Y$ . The strong contamination model is known to subsume the total variation contamination of Definition 1.2 (see, e.g., Section 2 of [DKK<sup>+</sup>16]). We note that our robust learning algorithm succeeds in this stronger contamination model, with the additional requirement that we can obtain two sets of independent  $\varepsilon$ -corrupted samples from  $\mathcal{M}$ . (The second set is needed to run a hypothesis testing routine after we obtain a small list of candidate hypotheses.)

Our algorithm works for any finite set of points in  $\mathbb{R}^d$  that satisfies a natural set of deterministic conditions. As we will show later in this section, these deterministic conditions are satisfied with high probability by a sufficiently large set of i.i.d. samples from any  $k$ -mixture of Gaussians.

**Condition 2.47** (Good Samples). Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians in  $\mathbb{R}^d$ . Let  $X$  be a set of  $n$  points in  $\mathbb{R}^d$ . We say that  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$  if there is a partition of  $X$  as  $X_1 \cup X_2 \cup \dots \cup X_k$  such that:

1. For all  $i \in [k]$  with  $w_i \geq \gamma$ , any positive integer  $m \leq t$ , and any  $v \in \mathbb{R}^d$ ,

$$\left| \frac{1}{n} \sum_{x \in X_i} \langle v, x - \mu_i \rangle^m - w_i \mathbb{E}_{x \sim \mathcal{N}(\mu_i, \Sigma_i)} [\langle v, x - \mu_i \rangle^m] \right| \leq w_i \gamma m! (v^T \Sigma_i v)^{m/2}.$$

2. For all  $i \in [k]$  and any halfspace  $H \subset \mathbb{R}^d$ , we have that  $|\frac{|X_i \cap H|}{n} - w_i \mathbb{P}_{x \sim \mathcal{N}(\mu_i, \Sigma_i)}[x \in H]| \leq \gamma$ .

We will also need the following consequences of Condition 2.47. The first one is immediate.

**Lemma 2.48.** Condition 2.47 is invariant under affine transformations. In particular, if  $A(x) : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  is an affine transformation, and if  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$ , then  $A(X)$  satisfies Condition 2.47 with respect to  $A(\mathcal{M})$  with parameters  $(\gamma, t)$ .

We note that the first part of Condition 2.47 implies that higher moment tensors are close in Frobenius distance.

**Lemma 2.49.** If  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  with parameters  $(\gamma, t)$ , then if  $w_i \geq \gamma$  for all  $i \in [k]$ , and if for some  $B \geq 0$  we have that  $\|\mu_i\|_2^2, \|\Sigma_i\|_{\text{op}} \leq B$  for all  $i \in [k]$ , then for all  $m \leq t$ , we have that:

$$\left\| \mathbb{E}_{x \in X} [x^{\otimes m}] - \mathbb{E}_{x \sim \mathcal{M}} [x^{\otimes m}] \right\|_F^2 \leq \gamma^2 m^{O(m)} B^m d^m.$$

We note that Condition 2.47 also behaves well with respect to taking submixtures.

**Lemma 2.50.** Let  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$ . Let  $S \subset [k]$  with  $\sum_{i \in S} w_i = w$ , and let  $\mathcal{M}' = \sum_{i \in S} (w_i/w) \mathcal{N}(\mu_i, \Sigma_i)$ . Then if  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$  for some  $\gamma < 1/(2k)$  with the corresponding partition being  $X = X_1 \cup X_2 \cup \dots \cup X_k$ , then  $X' = \bigcup_{i \in S} X_i$  satisfies Condition 2.47 with respect to  $\mathcal{M}'$  with parameters  $(O(k\gamma/w), t)$ .

Finally, we show that given sufficiently many i.i.d. samples from a  $k$ -mixture of Gaussians, Condition 2.47 holds with high probability.

**Lemma 2.51.** Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  and let  $n$  be an integer at least  $kt^C d^t / \gamma^3$ , for a sufficiently large universal constant  $C > 0$ , some  $\gamma > 0$ , and some  $t \in \mathbb{N}$ . If  $X$  consists of  $n$  i.i.d. samples from  $\mathcal{M}$ , then  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$  with high probability.

The proofs of the preceding lemmas can be found in Appendix A.



## 2.5 Hypothesis Selection

Our algorithm will require a procedure to select a hypothesis from a list of candidates that contains an accurate hypothesis. A number of such procedures are known in the literature (see, e.g., [Yat85, DL01, DDS12, DK14, DDS15, DKK<sup>+</sup>16]). Here we will use the following variant from [Kan20], showing that we can efficiently perform a hypothesis selection (tournament) step with access to  $\varepsilon$ -corrupted samples.

**Fact 2.52** (Robust Tournament, [Kan20]). *Let  $X$  be an unknown distribution,  $\eta \in (0, 1)$ , and let  $H_1, \dots, H_n$  be distributions with explicitly computable probability density functions that can be efficiently sampled from. Assume furthermore that  $\min_{1 \leq i \leq n} (d_{\text{TV}}(X, H_i)) \leq \eta$ . Then there exists an efficient algorithm that given access to  $O(\log(n)/\eta^2)$   $\varepsilon$ -corrupted samples from  $X$ , where  $\varepsilon \leq \eta$ , along with  $H_1, \dots, H_n$ , computes an  $m \in [n]$  such that with high probability we have that  $d_{\text{TV}}(X, H_m) = O(\eta)$ .*

## 3 List-Recovery of Parameters via Tensor Decomposition

In this section, we give an algorithm that takes samples from a  $k$ -mixture of Gaussians, whose component means and covariances are not too far from each other in natural norms, and outputs a dimension-independent size list of candidate  $k$ -tuples of parameters (i.e., means and covariances) one of which is guaranteed to be close to the true target  $k$ -tuple of parameters. Our approach involves a new tensor decomposition procedure that works in the absence of any non-degeneracy conditions on the components.

The goal of this section is to prove the following theorem:

**Theorem 3.1** (Recovering Candidate Parameters when Component Covariances are close in Frobenius Distance). *Fix any  $\alpha > \varepsilon > 0, \Delta > 0$ . There is an algorithm that takes input  $X$ , a sample from a  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  satisfying Condition 2.47 with parameters  $\gamma = \varepsilon d^{-8k} k^{-Ck}$ , for  $C$  a sufficiently large universal constant, and  $t = 8k$ , and let  $Y$  be an  $\varepsilon$ -corruption of  $X$ . If  $w_i \geq \alpha$ ,  $\|\mu_i\|_2 \leq \frac{2}{\sqrt{\alpha}}$  and  $\|\Sigma_i - I\|_F \leq \Delta$  for every  $i \in [k]$ , then, given  $k, Y$  and  $\varepsilon$ , the algorithm outputs a list  $L$  of at most  $\exp\left(\log(1/\varepsilon)(k + 1/\alpha + \Delta)^{O(k)} / \eta^2\right)$  candidate hypotheses (component means and covariances), such that with probability at least 0.99 there exist  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]} \subseteq L$  satisfying  $\|\mu_i - \hat{\mu}_i\|_2 \leq O\left(\frac{\Delta^{1/2}}{\alpha}\right) \eta^{G(k)}$  and  $\left\|\Sigma_i - \hat{\Sigma}_i\right\|_F \leq O(k^4) \frac{\Delta^{1/2}}{\alpha} \eta^{G(k)}$  for all  $i \in [k]$ . Here,  $\eta = (2k)^{4k} O(1/\alpha + \Delta)^{4k} \sqrt{\varepsilon}$ ,  $G(k) = \frac{1}{C^{k+1}(k+1)!}$ . The running time of the algorithm is  $\text{poly}(|L|, |Y|, d^k)$ .*

In the body of this section, we establish Theorem 3.1. The structure of this section is as follows: In Section 3.1, we describe our algorithm, which is then analyzed in Sections 3.2-3.6.

### 3.1 List-Decodable Tensor Decomposition Algorithm

In this section, we describe our tensor decomposition algorithm, which is given in pseudocode below (Algorithm 3.2).

**Algorithm 3.2** (List-Recovery of Candidate Parameters via Tensor Decomposition).

**Input:** An  $\varepsilon$ -corruption  $Y$  of a sample  $X$  from a  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$ .

**Requirements:** The guarantees of the algorithm hold if the mixture parameters and the sample  $X$  satisfy:

1.  $w_i \geq \alpha$  for all  $i \in [k]$ ,
2.  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$  for all  $i \in [k]$ ,
3.  $\|\Sigma_i - I\|_F \leq \Delta$  for all  $i \in [k]$ .
4.  $X$  satisfies Condition 2.47 with parameters  $(\gamma, t)$ , where  $\gamma = \varepsilon d^{-8k} k^{-Ck}$ , for  $C$  a sufficiently large universal constant, and  $t = 8k$ .

**Parameters:**  $\eta = (2k)^{4k} (Ck(1/\alpha + \Delta))^{4k} \sqrt{\varepsilon}$ ,  $D = C(k^4/(\alpha\sqrt{\eta}))$ ,  $\delta = 2\eta^{1/(C^{k+1}(k+1)!)}$ ,  $\ell' = 100 \log k (\eta/(k^5 (\Delta^4 + 1/\alpha^4)))^{-4k}$ , for some sufficiently large absolute constant  $C > 0$ ,  $\lambda = 4\eta$ ,  $\phi = 10(1 + \Delta^2)/(\sqrt{\eta}\alpha^5)$ .

**Output:** A list  $L$  of hypotheses such that there exists at least one,  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \leq k} \in L$ , satisfying:  
 $\|\mu_i - \hat{\mu}_i\|_2 \leq O\left(\frac{\Delta^{1/2}}{\alpha}\right) \eta^{G(k)}$  and  $\|\Sigma_i - \hat{\Sigma}_i\|_F \leq O(k^4) \frac{\Delta^{1/2}}{\alpha} \eta^{G(k)}$ , where  $G(k) = \frac{1}{C^{k+1}(k+1)!}$ .

**Operation:**

1. **Robust Estimation of Hermite Tensors:** For  $m \in [4k]$ , compute  $\hat{T}_m$  such that  $\max_{m \in [4k]} \|\hat{T}_m - \mathbb{E}[h_m(\mathcal{M})]\|_F \leq \eta$  using the robust mean estimation algorithm in Fact 2.36.
2. **Random Collapsing of Two Modes of  $\hat{T}_4$ :** Let  $L'$  be an empty list. Repeat  $\ell'$  times: For  $j \in [4k]$ , choose independent standard Gaussians in  $\mathbb{R}^d$ , denoted by  $x^{(j)}, y^{(j)} \sim \mathcal{N}(0, I)$ , and uniform draws  $a_1, a_2, \dots, a_t$  from  $[-D, D]$ . Let  $\hat{S}$  be a  $d \times d$  matrix such that for all  $r, s \in [d]$ ,  $\hat{S}(r, s) = \sum_{j \in [4k]} a_j \hat{T}_4(r, s, x^{(j)}, y^{(j)}) = \sum_{j \in [4k]} a_j \sum_{g, h \in [d]} \hat{T}_4(r, s, g, h) x^{(j)}(g) y^{(j)}(h)$ . Add  $\hat{S}$  to the list  $L'$ .
3. **Construct Low-Dimensional Subspace for Exhaustive Search:** Let  $V$  be the span of all singular vectors of the natural  $d \times d^{m-1}$  flattening of  $\hat{T}_m$  with singular values  $\geq \lambda$  for  $m \leq 4k$ . For each  $\hat{S} \in L'$ , let  $V'_\delta$  be the span of  $V$  plus all the singular vectors of  $\hat{S}$  with singular value larger than  $\delta^{1/4}$ .
4. **Enumerating Candidates in  $V'_\delta$ :** Initialize  $L$  to be the empty list. For each  $\hat{S} \in L'$ , let  $V_{\delta^{1/4}}$  be a  $\delta^{1/4}$ -cover of vectors in  $V'_\delta$  with  $\ell_2$ -norm at most  $2/\sqrt{\alpha}$ . Enumerate over vectors  $\hat{\mu}$  in  $V_{\delta^{1/4}}$ . Let  $k' = Ck^2$  and let  $C_{\delta^{1/4}}$  be a  $\delta^{1/4}$ -cover of the interval  $[-\phi, \phi]^{k'}$ . For  $\{\tau_j\}_{j \in [k']} \in C_{\delta^{1/4}}$  and for all  $\{v_j\}_{j \in [k']} \in V_{\delta^{1/4}}$ , let  $\hat{Q} = \sum_{j \in [k']} \tau_j v_j v_j^\top$ . Add  $\{\hat{\mu}, I + \hat{S} + \hat{Q}\}$  to  $L$ .



### 3.2 Analysis of Algorithm

We analyze the three main steps of Algorithm 3.2 in the following lemmas. We will prove the following three propositions in the subsequent subsections that analyze Steps 1, 2 and 3 of Algorithm 3.2. For Step 1, we show that when  $X$  satisfies Condition 2.47, the empirical estimates of the moment tensors obtained by applying the robust mean estimation algorithm to  $X$  are sufficiently close to the moment tensors of the input mixture  $\mathcal{M}$ .

**Proposition 3.3** (Robustly Estimating Hermite Polynomial Tensors). *For any integer  $m \leq 4k$ , and  $\Delta \in \mathbb{R}_+$ , there exists an algorithm with running time  $\text{poly}_m(d/\varepsilon)$  that takes an  $\varepsilon$ -corruption  $Y$  of  $X$ , a set satisfying Condition 2.47 with respect to  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  with parameters  $\gamma = \varepsilon d^{-m} m^{-Cm}$ , for  $C$  a sufficiently large constant, and  $t = 2m$ . If  $w_i \geq \alpha$ ,  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$ , and  $\|\Sigma_i - I\|_F \leq \Delta$  for each  $i \in [k]$ , then the algorithm outputs a tensor  $\hat{T}_m$  such that  $\left\| \hat{T}_m - \mathbb{E}[h_m(\mathcal{M})] \right\|_F \leq \eta$ , for  $\eta = \mathcal{O}(m(1 + 1/\alpha + \Delta))^m \sqrt{\varepsilon}$ .*

The proof of Proposition 3.3 is deferred to Section 3.3.

Next, we analyze Step 2 of the algorithm and prove that, with non-negligible probability, randomly collapsing two modes of  $\hat{T}_4$  yields a matrix  $\hat{S}$  such that  $\hat{S} - (\Sigma_i - I) = P_i + Q_i$ , where  $P_i$  has small Frobenius norm and  $Q_i$  is a rank- $\mathcal{O}(k^2)$  matrix.

**Proposition 3.4** (Tensor Decomposition up to Low-Rank Error). *Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians satisfying  $w_i \geq \alpha$ ,  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$ , and  $\|\Sigma_i - I\|_F \leq \Delta$  for each  $i \in [k]$ . For  $0 < \eta < 1$ , let  $\hat{T}_4$  be a tensor such that  $\left\| \mathbb{E}[h_4(\mathcal{M})] - \hat{T}_4 \right\|_F \leq \eta$ , and let  $D$  be a sufficiently large constant multiple of  $k^4/(\alpha\sqrt{\eta})$ . For all  $j \in [4k]$ , let  $x^{(j)}, y^{(j)} \sim \mathcal{N}(0, I)$  be independent and  $a_j \sim \mathcal{U}[-D, D]$ , where  $\mathcal{U}[-D, D]$  is the uniform distribution over the interval  $[-D, D]$ , and let  $\hat{S} = \sum_{j \in [4k]} a_j \hat{T}_4(\cdot, \cdot, x^{(j)}, y^{(j)})$ . Then, for each  $i \in [k]$ , with probability at least  $(\eta/(k^5(\Delta^4 + 1/\alpha^4)))^{4k}$ , over the choice of  $x^{(j)}, y^{(j)}$  and  $a_j$ , we have that  $\hat{S} - (\Sigma_i - I) = P_i + Q_i$ , where  $\|P_i\|_F = \mathcal{O}(\sqrt{\eta/\alpha})$ ,  $\|Q_i\|_F = \mathcal{O}\left(\frac{1+\Delta^2}{\sqrt{\eta\alpha^3}}\right)$  and  $\text{rank}(Q_i) = \mathcal{O}(k^2)$ .*

The proof of Proposition 3.4 is given in Section 3.4.

Finally, in Step 3, for any  $\hat{S}$  such that  $\hat{S} - (\Sigma_i - I) = P_i + Q_i$ , where  $P_i$  has small Frobenius norm and  $Q_i$  is a rank  $\mathcal{O}(k^2)$  matrix, we find a low-dimensional subspace  $V'$  such that the range space of  $Q_i$  is approximately contained in  $V'$ . We will use  $V'$  to exhaustively search for  $\mathcal{O}(k^2)$  rank matrices to find candidates for  $Q_i$ .

**Proposition 3.5** (Low-Dimensional Subspace  $V'$  for Exhaustive Search). *Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians satisfying  $w_i \geq \alpha$ ,  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$ , and  $\|\Sigma_i - I\|_F \leq \Delta$  for each  $i \in [k]$ . Let  $\left\| \hat{T}_m - \mathbb{E}[h_m(\mathcal{M})] \right\|_F \leq \eta$ , for each  $1 \leq m \leq 4k$ , and some  $\eta > 0$ . Let  $V$  be the span of all the left singular vectors of the  $d \times d^{m-1}$  matrix obtained by the natural flattening of  $\hat{T}_m$  with singular values at least  $2\eta$ . For each  $1 \leq i \leq k$ , let  $S_i = \Sigma_i - I$  and  $\hat{S}_i$  be a  $d \times d$  matrix such that  $\hat{S}_i - S_i = P_i + Q_i$ , where  $\|P_i\|_F \leq \mathcal{O}(\sqrt{\eta/\alpha})$ ,  $Q_i$  has rank  $\mathcal{O}(k^2)$ , and  $\|Q_i\|_F \leq \mathcal{O}\left(\frac{1+\Delta^2}{\sqrt{\eta\alpha^3}}\right)$ . Let  $V'$  be the span of  $V$  plus all singular vectors of  $\hat{S}_i$  of singular values at least  $\delta$  for all  $i$ . Then, for  $\delta = 2\eta^{1/(C^{k+1}(k+1)!)}$  with a sufficiently large constant  $C > 0$ , we have that:*

1.  $\dim(V') \leq \left( \mathcal{O}(k(1 + 1/\alpha + \Delta))^{4k+5} \right) / \eta^2$ .

2. There is a vector  $\mu'_i \in V'$  such that  $\|\mu_i - \mu'_i\|_2^2 \leq \frac{20}{\alpha^2} \sqrt{\delta} \Delta$ .
3. There are  $q = O(k^2)$  unit vectors  $v_1, v_2, \dots, v_q \in V'$  and scalars  $\tau_1, \tau_2, \dots, \tau_q \in [-10(1 + \Delta^2)/(\sqrt{\eta}\alpha^5), 10(1 + \Delta^2)/(\sqrt{\eta}\alpha^5)]$  such that  $\|Q_i - \sum_{i=1}^q \tau_i v_i v_i^\top\|_F \leq O\left(\frac{k^2}{\alpha} \delta^{1/4} \Delta^{1/2}\right)$ .

The proof of Proposition 3.5 is given in Section 3.5.

We can now use these propositions to complete the proof of Theorem 3.1.

*Proof of Theorem 3.1.* Using Proposition 3.3, Step 1 of the algorithm outputs estimates  $\hat{T}_i$  for  $i \in [4k]$  such that  $\max_{m \in [4k]} \|\hat{T}_m - \mathbb{E} h_m(\mathcal{M})\|_F \leq \eta$ . Next, by the standard coupon collector analysis, using Proposition 3.4 and repeating Step 2 of the algorithm  $\ell' = 100 \log k (\eta / (k^5 (\Delta^4 + 1/\alpha^4)))^{-4k}$  times, guarantees that with probability at least  $1 - 1/(100k)^{100}$ , for every  $1 \leq i \leq k$ , there are  $\hat{S}_i \in L$  such that  $\hat{S}_i - (\Sigma_i - I) = P_i + Q_i$  for  $P_i, Q_i$  satisfying  $\|P_i\|_F \leq \sqrt{\eta/\alpha}$ ,  $\|Q_i\|_F \leq \frac{1+\Delta^2}{\sqrt{\eta}\alpha^5}$  and  $Q_i$  has rank  $O(k^2)$ .

Next, Proposition 3.5 implies that for every such  $\hat{S}_i \in L'$ , we can construct a subspace  $V' = V'_{\hat{S}_i}$  of dimension  $O((k(1 + 1/\alpha + \Delta))^{4k+5}/\eta^2)$  such that  $V'$  contains  $\mu'_i$  that satisfies  $\|\mu_i - \mu'_i\|_2^2 \leq \frac{\Delta}{\alpha^2} \cdot \sqrt{\delta}$ , and there is a rank  $O(k^2)$  matrix  $\hat{Q}_i$  with range space contained in  $V'$  such that  $\|Q_i - \hat{Q}_i\|_F \leq O\left(\frac{k^2}{\alpha} \delta^{1/4} \Delta^{1/2}\right)$ .

Now, let  $V_\tau \subseteq V'$  be a  $\tau = \delta^{1/4}$ -cover, in  $\ell_2$ -norm, of vectors with  $\ell_2$  norm at most  $2/\sqrt{\alpha}$  in  $V'$ . Then, since  $\|\mu_i\|_2 \leq \frac{2}{\sqrt{\alpha}}$ , there is a vector  $\hat{\mu}_i \in V_\tau$  such that  $\|\mu_i - \hat{\mu}_i\|_2^2 \leq \tau + \frac{20}{\alpha^2} \sqrt{\delta} \Delta \leq \frac{40}{\alpha^2} \sqrt{\delta} \Delta$ .

Further, there exist  $\tau_1, \tau_2, \dots, \tau_{O(k^2)}$  in a  $\tau$ -cover of  $[-10(1 + \Delta^2)/(\sqrt{\eta}\alpha^5), 10(1 + \Delta^2)/(\sqrt{\eta}\alpha^5)]$  and vectors  $v_1, v_2, \dots, v_{O(k^2)} \in V_\tau$  such that  $\left\| \sum_{i=1}^{O(k^2)} \tau_i v_i v_i^\top - Q_i \right\|_F \leq O(k^4 \delta^{1/4} \Delta^{1/2} / \alpha)$ . In particular,  $\hat{\Sigma}_i = I + \hat{S}_i - \sum_{i=1}^{O(k^2)} \tau_i v_i v_i^\top$  satisfies

$$\|\hat{\Sigma}_i - \Sigma_i\|_F = O(\sqrt{\eta}) + O\left(\frac{k^4 \delta^{1/4} \Delta^{1/2}}{\alpha}\right) = O\left(\frac{k^4 \delta^{1/4} \Delta^{1/2}}{\alpha}\right). \quad (3.1)$$

The size of this search space for every fixed  $\hat{S} \in L'$  can be bounded above by  $\left(\frac{1+\Delta^2}{\delta\alpha^5}\right)^{O(k^5 \dim(V'))}$ . Thus, the size of  $L$  can be bounded from above by

$$k^5 \left(\frac{\Delta^4}{\eta} + \frac{1}{\alpha^4 \eta}\right)^{4k} \cdot \left(\frac{1 + \Delta^2}{\delta \sqrt{\eta} \alpha^5}\right)^{O(k^5 \dim(V'))} \leq \exp\left(\log(1/\varepsilon) (k + 1/\alpha + \Delta)^{O(k)} / \eta^2\right).$$

This completes the proof.  $\square$

### 3.3 Robust Estimation of Hermite Tensors

In this section, we will prove Proposition 3.3.

*Proof of Proposition 3.3.* Consider the uniform distribution on the uncorrupted sample  $X$ . We want to analyze the effect of applying the robust mean estimation algorithm (Fact 2.36) to the points  $h_m(x)$ , for  $x \in X$ . In order for us to apply Fact 2.36, we need to ensure that the uniform distribution

on  $\{h_m(x)\}_{x \in X}$  has bounded covariance. This step gives us a good approximation to  $\mathbb{E}_{x \sim_u X} h_m(x)$ . In order for us to obtain an approximation to  $\mathbb{E} h_m(\mathcal{M})$ , we need to bound the difference between  $\mathbb{E} h_m(\mathcal{M})$  and  $\mathbb{E}_{x \sim_u X} h_m(x)$ . We will do both these steps below.

The second part is immediate. By the definition of  $h_m(X)$ , we have that

$$\left\| \frac{1}{|X|} \sum_{x \in X} h_m(x) - \mathbb{E} h_m(\mathcal{M}) \right\|_F \leq \sum_{j \leq m/2} m^{2j} d^j \left\| \frac{1}{|X|} \sum_{x \in X} x^{\otimes(m-2j)} - \mathbb{E} \mathcal{M}^{\otimes(m-2j)} \right\|_F.$$

By Lemma 2.49, this is at most  $O(1 + \Delta + 1/\alpha)^m m^{O(m)} d^{m/2} \gamma \leq \eta/2$ . We note that a similar argument bounds

$$\left\| \frac{1}{|X|} \sum_{x \in X} h_m(x) \otimes h_m(x) - \mathbb{E} h_m(\mathcal{M}) \otimes h_m(\mathcal{M}) \right\|_F \leq \eta^2.$$

Let us now verify the first part. We proceed via bounding the operator norm of the covariance of  $h_m(\mathcal{M})$ . We can then use the bound on the Frobenius norm  $\left\| \frac{1}{|X|} \sum_{x \in X} h_m(x) \otimes h_m(x) - \mathbb{E} h_m(\mathcal{M}) \otimes h_m(\mathcal{M}) \right\|_F$  to get a bound on  $\left\| \frac{1}{|X|} \sum_{x \in X} h_m(x) h_m(x)^\top \right\|_{\text{op}}$  (the operator norm of the canonical square flattening of the of the  $2m$ -th empirical Hermite moment tensor of  $X$ ). This will complete the proof.

Let  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be the components of  $\mathcal{M}$ . We have that

$$\begin{aligned} \text{Cov}(h_m(\mathcal{M})) &= \sum_{i \in [k]} w_i \text{Cov}(h_m(G_i)) \\ &+ \frac{1}{2} \sum_{i, j \in [k]} w_i w_j (\mathbb{E}[h_m(G_i)] - \mathbb{E}[h_m(G_j)]) (\mathbb{E}[h_m(G_i)] - \mathbb{E}[h_m(G_j)])^\top. \end{aligned} \quad (3.2)$$

By Lemma 2.8, we have that for all  $i \in [k]$ , it holds

$$\|\text{Cov}(h_m(G_i))\|_{\text{op}} = O\left(m(1 + \|\mu_i\|_2 + \|\Sigma_i - I\|_F)\right)^{2m} = O\left(m(1 + 2/\sqrt{\alpha} + \Delta)\right)^{2m},$$

where for any matrix  $M$ ,  $\|M\|_{\text{op}} = \max_{\|u\|_2=1} \|Mu\|_2$  is the operator norm of the matrix. Further, for any  $i, j \in [k]$ ,

$$\begin{aligned} \left\| (\mathbb{E}[h_m(G_i)] - \mathbb{E}[h_m(G_j)]) (\mathbb{E}[h_m(G_i)] - \mathbb{E}[h_m(G_j)])^\top \right\|_{\text{op}} &= \|\mathbb{E}[h_m(G_i)] - \mathbb{E}[h_m(G_j)]\|_2 \\ &= O(m(1 + 1/\alpha + \Delta))^{2m}. \end{aligned} \quad (3.3)$$

This claim follows from the triangle inequality of the operator norm.  $\square$

### 3.4 List-Recovery of Covariances up to Low-Rank Error

In this section, we prove Proposition 3.4. We first set some useful notation. We will write  $S_i \stackrel{\text{def}}{=} \Sigma_i - I$  throughout this section. We will also use  $S'_i$  to denote  $S_i + \mu_i \otimes \mu_i$ .

We first show that for every  $i$ , there exists a matrix  $P$  such that  $\left(\sum_{i \in [k]} w_i S'_i \otimes S'_i\right)(\cdot, \cdot, P)$  is close to  $S'_i$ .

**Lemma 3.6** (Existence of a 2-Tensor). *Under the hypothesis of Proposition 3.4, for each  $i \in [k]$ , there exists a matrix  $P$  such that  $\|P\|_F = O(1/(\sqrt{\eta}\alpha))$  and  $\|T'_4(\cdot, \cdot, P) - S'_i\|_F = O(\sqrt{\eta/\alpha})$ , where  $T'_4 = (\sum_{i \in [k]} w_i S'_i \otimes S'_i)$ .*

Note that throughout this section it will be useful to think of  $T'_4$  as a  $d^2 \times d^2$  matrix rather than as a tensor. In this case, we can think of  $T'_4$  as  $\sum_{i=1}^k w_i (S'_i)(S'_i)^T$ . From standard facts about positive semidefinite matrices it follows that  $S'_i$  is in the image of  $T'_4$ , and Lemma 3.6 is just a slightly robustified version of this (saying that we can find an approximate preimage that is not itself too large).

The proof of this Lemma 3.6 will involve linear programming duality with an infinite system of constraints. As the application of duality with infinitely many constraints has some technical issues, we state below an appropriate version of duality.

**Fact 3.7** (Linear Programming Duality for Compact, Convex Constraint Sets). *Let  $K \subset \mathbb{R}^{n+1}$  be a compact convex set. There exists an  $x \in \mathbb{R}^n$  so that  $(x, 1) \cdot z > 0$  for all  $z \in K$  if and only if there is no element  $(0, 0, \dots, 0, a) \in K$  for any  $a \leq 0$ .*

This fact can be proved by noting that if no such  $a$  exists, there must be a hyperplane separating  $K$  from the set of such points  $(0, a)$ . This separating hyperplane will be of the form  $(z, y) \in H$  if and only if  $y = x \cdot z$  for some  $x$  and this  $x$  will provide the solution to the linear system.

*Proof of Lemma 3.6.* To show that such a  $P$  exists for each  $i$ , we apply linear programming duality. In particular, the conditions imposed on  $P$  define a linear program, which has a feasible solution unless there is a solution to the dual linear program. For sufficiently large constants  $c_1$  and  $c_2$ , consider the following primal in the variable  $P$ :

$$\langle v, P \rangle \leq \frac{c_1}{\sqrt{\eta}\alpha} \|v\|_F \quad \forall v \in \mathbb{R}^{d \times d} \quad (3.4)$$

$$\langle u, T'_4(\cdot, \cdot, P) - S'_i \rangle \leq c_2 \sqrt{\eta} \|u\|_F \quad \forall u \in \mathbb{R}^{d \times d}. \quad (3.5)$$

It is not hard to see that  $\|P\|_F \leq \frac{c_1}{\sqrt{\eta}\alpha}$  if and only if (3.4) holds for all  $v$  and  $\|T'_4(\cdot, \cdot, P) - S'_i\|_F \leq c_2 \sqrt{\eta/\alpha}$  if and only if (3.5) holds for all  $u$ . Throughout the proof, we suggest that the reader think of  $u$  and  $v$  as vectors in  $d^2$ -dimensional vector space.

Our goal is to show that there exists a feasible solution  $P$  such that (3.4) and (3.5) hold simultaneously for all  $u, v \in \mathbb{R}^{d \times d}$ . We first note that this is equivalent to saying that

$$\langle v, P \rangle + \langle u, T'_4(\cdot, \cdot, P) \rangle - \langle u, S'_i \rangle \leq \frac{c_1}{\sqrt{\eta}\alpha} \|v\|_F + c_2 \sqrt{\eta} \|u\|_F, \quad (3.6)$$

for all  $u, v \in \mathbb{R}^{d \times d}$ . This is not quite in the form necessary to apply Fact 3.7, so we note that this is in turn equivalent to saying that

$$\langle v, P \rangle + \langle u, T'_4(\cdot, \cdot, P) \rangle - \langle u, S'_i \rangle \leq 1, \quad (3.7)$$

for all  $u, v \in \mathbb{R}^{d \times d}$  so that  $\frac{c_1}{\sqrt{\eta}\alpha} \|v\|_F + c_2 \sqrt{\eta} \|u\|_F \leq 1$ , and  $u \in \text{span}\{S'_i\}$ . As this is a convex set of linear equations, we have by Fact 3.7 that there exists such a  $P$  unless there exists such a pair of

$u$  and  $v$  so that the coefficient of  $P$  in Equation (3.7) is 0 and so that the resulting inequality of constants is either false or holds with equality. In particular, the coefficient of  $P$  vanishes if and only if  $v = -T'_4(u, \cdot, \cdot)$ . We then get a contradiction only if for some  $u \in \text{span}\{S'_i\}$

$$-\langle u, S'_i \rangle \geq 1 \geq \frac{c_1}{\sqrt{\eta}\alpha} \|T'_4(u, \cdot, \cdot)\|_F + c_2\sqrt{\eta} \|u\|_F. \quad (3.8)$$

We claim that this is impossible.

In particular, squaring Equation (3.8) would give

$$\begin{aligned} \langle u, S'_i \rangle^2 &\geq \left( \frac{c_1}{\sqrt{\eta}\alpha} \|T'_4(u, \cdot, \cdot)\|_F + c_2\sqrt{\eta} \|u\|_F \right)^2 \\ &\geq \frac{c}{\alpha} \|T'_4(u, \cdot, \cdot)\|_F \cdot \|u\|_F, \end{aligned} \quad (3.9)$$

for some large enough constant  $c > 1$ , where the last inequality follows from the AM-GM inequality. However, using the dual characterization of the Frobenius norm, we have

$$\|T'_4(u, \cdot, \cdot)\|_F \geq \frac{\langle u, T'_4(u, \cdot, \cdot) \rangle}{\|u\|_F} \geq \frac{w_i}{\|u\|_F} \langle u, S'_i \rangle^2, \quad (3.10)$$

where the last inequality follows from  $T'_4$  containing a  $w_i S_i \otimes S_i$  term, and the other terms contributing non-negatively. Rearranging Equation (3.10), we have

$$\langle u, S'_i \rangle^2 \leq \frac{1}{w_i} \|T'_4(u, \cdot, \cdot)\|_F \|u\|_F \leq \frac{1}{\alpha} \|T'_4(u, \cdot, \cdot)\|_F \|u\|_F.$$

This contradicts Equation (3.9) unless  $T'_4(u, \cdot, \cdot) = 0$ . This therefore suffices to prove the feasibility of the primal.  $\square$

We have thus shown that there is some matrix  $P$  so that  $T'_4(P, \cdot, \cdot)$  suffices for our purposes. We need to show that our appropriate random linear combination of  $x^{(j)} \otimes y^{(j)}$  suffices. In fact, we will show that with reasonably high probability over our choice of  $x^{(j)}, y^{(j)}$  that there is some linear combination of the  $x^{(j)} \otimes y^{(j)}$  (with coefficients that are not too large) so that their projection onto the space spanned by the  $S'_i$  (which is all that matters when applying  $T'_4$ ) equal to  $P$ .

For the sake of intuition, we note that if we removed the bound on the coefficients, we would need that the projections of the  $x^{(j)} \otimes y^{(j)}$  spanned  $\text{span}\{S'_i\}$ . Since there are at least  $k$  of them, this will hold unless there is some  $v \in \text{span}\{S'_i\}$  so that  $v$  is orthogonal to all of the  $x^{(j)} \otimes y^{(j)}$ . This shouldn't happen because each  $x^{(j)} \otimes y^{(j)}$  is very unlikely to be orthogonal to  $v$ .

To deal with the constraint that the coefficients are not too large, we use linear programming duality to show that there will be a solution unless there is some  $v$  that is *nearly* orthogonal to all of the  $x^{(j)} \otimes y^{(j)}$ . Again, this is unlikely to happen for any individual term, and thus, by independence, highly unlikely to happen for all  $j$  simultaneously. Combining this with a cover argument will give our proof.

**Lemma 3.8** (Existence of a Bi-Linear Form). *Given the preconditions in Proposition 3.4, with probability at least 99/100 over the choice of  $x^{(j)}, y^{(j)}$ , there exist  $b_j \in [-D, D]$  for  $j \in [4k]$ , where  $D = \mathcal{O}(k^4/(\sqrt{\eta}\alpha))$ , such that the projection of  $\sum_{j=1}^t b_j x^{(j)} \otimes y^{(j)}$  onto the space spanned by the  $S'_i$  is  $P$ , where  $P$  satisfies the conclusion of Proposition 3.6.*

*Proof.* To prove this lemma, we again use a linear programming based argument. Consider the following (primal) linear program in the variables  $b_j$ , for  $j \in [4k]$ :

$$\sum_{j \in [4k]} b_j \langle S'_i, x^{(j)} \otimes y^{(j)} \rangle = \langle S'_i, P \rangle \quad \forall i \in [k] \quad (3.11)$$

$$-D \leq b_j \leq D \quad \forall j \in [4k] \quad (3.12)$$

We note that a set of  $b_j$  satisfying Equation (3.11) will have the projection of  $\sum_{j \in [4k]} b_j x^{(j)} \otimes y^{(j)}$  onto the span of the  $S'_i$  be the same as the projection of  $P$ , and that if the  $b_j$ 's satisfy Equation (3.12) then we will have  $|b_j| \leq D$  for all  $j$ . Thus, it suffices to show that with high probability over our choice of  $x^{(j)}$  and  $y^{(j)}$  that the above system is feasible.

We will show this by linear programming duality (since this is now a finite system of equations, we can use standard results rather than Fact 3.7). In particular, we have that Equations (3.11) and (3.12) are simultaneously satisfiable unless there are real numbers  $c_i$  and non-negative real numbers  $z_j, z'_j$  so that

$$\sum_{i=1}^k c_i \sum_{j \in [4k]} b_j \langle S'_i, x^{(j)} \otimes y^{(j)} \rangle + \sum_{j \in [4k]} (z_j - z'_j) b_j \leq \sum_{i=1}^k c_i \langle S'_i, P \rangle + \sum_{j \in [4k]} (z_j + z'_j) D$$

yields a contradiction. Setting  $v = \sum_{i=1}^k c_i S'_i$ , the above simplifies to

$$\sum_{j \in [4k]} b_j \left( \langle v, x^{(j)} \otimes y^{(j)} \rangle + z_j - z'_j \right) \leq \langle v, P \rangle + \sum_{j \in [4k]} (z_j + z'_j) D \quad (3.13)$$

We note that in order for Equation (3.13) to be a contradiction, it must be the case that the coefficients of  $b_j$  are all 0. In particular, we must have

$$z'_j - z_j = \langle v, x^{(j)} \otimes y^{(j)} \rangle$$

for all  $j$ . In particular, this means that

$$z_j + z'_j \geq |\langle v, x^{(j)} \otimes y^{(j)} \rangle|.$$

In such a case, the right hand side of Equation (3.13) will be at least

$$\langle v, P \rangle + \sum_{j \in [4k]} |\langle v, x^{(j)} \otimes y^{(j)} \rangle| D$$

Therefore, Equation (3.13) can only yield a contradiction if there exists a  $v \in \text{span}\{S'_i\}$  so that

$$\langle v, P \rangle < - \sum_{j \in [4k]} |\langle v, x^{(j)} \otimes y^{(j)} \rangle| D. \quad (3.14)$$

We want to show that with high probability over our choice of  $x^{(j)}, y^{(j)}$  that there is no  $v \in \text{span}\{S'_i\}$  satisfying Equation (3.14). In fact, we will show that for every such  $v$  that

$$\sum_{j \in [4k]} |\langle v, x^{(j)} \otimes y^{(j)} \rangle| \geq \frac{c_1}{\sqrt{\eta} \alpha} \|v\|_F.$$

We can scale  $v$  so that  $\|v\|_F = 1$ , and it suffices to show that

$$\sum_{j \in [4k]} |\langle \tilde{v}, x^{(j)} \otimes y^{(j)} \rangle| \geq \left( \frac{c_1}{\sqrt{\eta} \alpha D} \right) \quad (3.15)$$

holds for all unit vectors  $v$  in  $\text{span}\{S'_i\}$  with high probability.

Since we need to show that infinitely many equations all hold with high probability, we will use a cover argument. In particular, we can construct  $\mathcal{C}$ , a  $\tau$ -cover for all unit vectors  $v$  in the span of the  $S'_i$ , where we take  $\tau = \left( \frac{c'_1}{k^2 \sqrt{\eta} \alpha D} \right)$ . Since this is a cover of a unit sphere in a  $k$ -dimensional subspace, we can construct such a cover so that  $|\mathcal{C}| = \mathcal{O}(1/\tau)^k$ . Replacing  $v$  with the closest point in  $\mathcal{C}$ , denoted by  $v'$ , it suffices to show that with high probability for all  $v$  that

$$\sum_{j \in [4k]} |\langle v, x^{(j)} \otimes y^{(j)} \rangle| \geq \sum_{j \in [4k]} |\langle v', x^{(j)} \otimes y^{(j)} \rangle| - \sum_{j \in [4k]} |\langle v - v', x^{(j)} \otimes y^{(j)} \rangle| \geq \left( \frac{2c_1}{\sqrt{\eta} \alpha D} \right). \quad (3.16)$$

We begin by bounding the terms

$$\sum_{j \in [4k]} |\langle v - v', x^{(j)} \otimes y^{(j)} \rangle|.$$

For this we notice by Cauchy-Schwartz that each term is at most  $\|v - v'\|_F$  times the Frobenius norm of the projection of  $x^{(j)} \otimes y^{(j)}$  onto the span of the  $S'_i$ . We note that for any  $k$ -dimensional subspace  $W$  with orthonormal basis  $w_1, \dots, w_k$  we have that

$$\begin{aligned} \mathbb{E} \left[ \left\| \text{Proj}_W(x^{(j)} \otimes y^{(j)}) \right\|_F^2 \right] &= \sum_{i=1}^k |\langle w_i, x^{(j)} \otimes y^{(j)} \rangle|^2 \\ &= k. \end{aligned}$$

Therefore, with high probability over the choice of  $x^{(j)}, y^{(j)}$  each of the projections of  $x^{(j)} \otimes y^{(j)}$  onto the span of the  $S'_i$  has Frobenius norm  $\tilde{\mathcal{O}}(\sqrt{k})$ . Therefore, if this condition holds over our choice of  $x^{(j)}$  and  $y^{(j)}$ , we can show Equation (3.16) if we can show that

$$\sum_{j \in [4k]} |\langle v', x^{(j)} \otimes y^{(j)} \rangle| \geq \left( \frac{c_1}{\sqrt{\eta} \alpha D} \right) \geq \left( \frac{2c_1}{\sqrt{\eta} \alpha D} \right) - \tau \tilde{\mathcal{O}}(k^{3/2}) \quad (3.17)$$

for all  $v' \in \mathcal{C}$ .

Each term in  $\sum_{j \in [4k]} \langle v', x^{(j)} \otimes y^{(j)} \rangle$  is a random bi-linear form given by  $z_j = \sum_{\ell, p \in [d]} v'_{\ell, p} x_{\ell}^{(j)} y_p^{(j)}$ . Then, we have that  $\mathbb{E}[z_j] = 0$  and

$$\begin{aligned} \mathbb{E}[z_j^2] &= \mathbb{E} \left[ \left( \sum_{\ell, p \in [d]} v'_{\ell, p} x_{\ell}^{(j)} y_p^{(j)} \right)^2 \right] = \sum_{\ell, \ell', p, p'} \mathbb{E} \left[ v'_{\ell, p} v'_{\ell', p'} x_{\ell}^{(j)} x_{\ell'}^{(j)} y_p^{(j)} y_{p'}^{(j)} \right] \\ &= \sum_{\ell, p \in [d]} (v'_{\ell, p})^2 \cdot \mathbb{E} \left[ (x_{\ell}^{(j)})^2 \right] \cdot \mathbb{E} \left[ (y_p^{(j)})^2 \right] \\ &= 1, \end{aligned}$$



where the last equality follows from  $\tilde{v}'_F = 1$ .

Using Lemma 2.10 with  $\zeta = \frac{2c_1}{\sqrt{\eta}\alpha D}$ ,

$$\mathbb{P} \left[ |z_j| \leq \frac{c_1}{\sqrt{\eta}\alpha D} \right] \leq c_5 \left( \frac{2c_1}{\sqrt{\eta}\alpha D} \right)^{1/2}. \quad (3.18)$$

However, we note that Equation (3.17) will hold unless  $|z_j| \leq \frac{c_1}{\sqrt{\eta}\alpha D}$  for all  $j \in [4k]$ . Since the  $z_j$ 's are independent, we conclude that

$$\mathbb{P} \left[ \sum_{j \in [4k]} |\langle v', x^{(j)} \otimes y^{(j)} \rangle| \leq \frac{c_1}{\sqrt{\eta}\alpha D} \right] \leq O \left( \frac{c_1}{\sqrt{\eta}\alpha D} \right)^{2k}. \quad (3.19)$$

Since the above argument holds for any  $v' \in \mathcal{C}$ , we can union bound over all elements in the cover  $\mathcal{C}$ , and the probability that there exists a  $\tilde{v}'$  in the cover that does not satisfy Equation (3.17) is at most  $O(k^2 \sqrt{\eta}\alpha D)^k \cdot O \left( \frac{c_1}{\sqrt{\eta}\alpha D} \right)^{2k}$ . Setting  $D$  to be a sufficiently large multiple of  $(k^4 / (\sqrt{\eta}\alpha))$  suffices to conclude that with probability at least  $1 - 1/\text{poly}(k)$ , the primal is feasible.  $\square$

*Proof of Proposition 3.4.* We begin by bounding the Frobenius norm of  $\hat{T}_4$ . Let  $T_4 = \mathbb{E}[h_4(X)]$ . It then follows from Lemma 2.6 that

$$T_4 = \text{Sym} \left( \sum_{i=1}^k w_i (3S_i \otimes S_i + 6S_i \otimes \mu_i^{\otimes 2} + \mu_i^{\otimes 4}) \right). \quad (3.20)$$

Further,  $\|S_i \otimes S_i\|_F \leq \|S_i\|_F^2 \leq \Delta^2$ ,  $\|S_i \otimes \mu_i^{\otimes 2}\|_F \leq \|S_i\|_F \|\mu_i\|_2^2 \leq 4\Delta/\alpha$ , and  $\|\mu_i^{\otimes 4}\|_F \leq \|\mu_i\|_2^4 \leq 16/\alpha^2$ . Since  $T_4$  is an average of terms of the form  $S_i^{\otimes 2}$ ,  $S_i \otimes \mu_i^{\otimes 2}$  and  $\mu_i^{\otimes 4}$ , and each such term is upper bounded, we can conclude that  $\|T_4\|_F = O(\Delta^2 + 1/\alpha^2)$ , and by the triangle inequality that  $\|\hat{T}_4\|_F \leq O(\Delta^2 + 1/\alpha^2 + \eta)$ . Let  $S'_i = S_i + \mu_i^{\otimes 2}$  and let  $T'_4 := \sum_{i=1}^k w_i (S'_i \otimes S'_i)$ . We can then rewrite Equation (3.20) as follows:

$$T_4 = \text{Sym} \left( \sum_{i=1}^k w_i (3S'_i \otimes S'_i - 2\mu_i^{\otimes 4}) \right). \quad (3.21)$$

For  $j \in [4k]$ , let  $x^{(j)}, y^{(j)} \sim \mathcal{N}(0, I)$ . Collapsing two modes of  $\hat{T}_4$ , it follows from Equation (3.21) that for any fixed  $j$ ,

$$\begin{aligned} \hat{T}_4(\cdot, \cdot, x^{(j)}, y^{(j)}) &= (\hat{T}_4 - T_4)(\cdot, \cdot, x^{(j)}, y^{(j)}) + T_4(\cdot, \cdot, x^{(j)}, y^{(j)}) \\ &= (\hat{T}_4 - T_4)(\cdot, \cdot, x^{(j)}, y^{(j)}) + \text{Sym} \left( \sum_{i=1}^k w_i (3S'_i \otimes S'_i - 2\mu_i^{\otimes 4}) \right) (\cdot, \cdot, x^{(j)}, y^{(j)}) \\ &= (\hat{T}_4 - T_4 + T'_4)(\cdot, \cdot, x^{(j)}, y^{(j)}) + \sum_{i \in [k]} w_i (S'_i x^{(j)}) \otimes (S'_i y^{(j)}) \\ &\quad + \sum_{i \in [k]} w_i (S'_i y^{(j)}) \otimes (S'_i x^{(j)}) + \sum_{i \in [k]} w_i (-2\mu_i^{\otimes 2} \langle \mu_i, x^{(j)} \rangle \langle \mu_i, y^{(j)} \rangle), \end{aligned} \quad (3.22)$$



where we use that  $\text{Sym}(\cdot)$  is a linear operator satisfying  $\text{Sym}(\mu_i^{\otimes 4}) = \mu_i^{\otimes 4}$ , and

$$\text{Sym}(S'_i \otimes S'_i) = \frac{1}{3}S'_i \otimes S'_i + \frac{1}{3}S'_i \oplus S'_i + \frac{1}{3}S'_i \ominus S'_i$$

where for indices  $(i_1, i_2, i_3, i_4)$ ,  $(S'_i \oplus S'_i)(i_1, i_2, i_3, i_4) = (S'_i \otimes S'_i)(i_1, i_3, i_2, i_4)$  and  $(S'_i \ominus S'_i)(i_1, i_2, i_3, i_4) = (S'_i \otimes S'_i)(i_1, i_4, i_2, i_3)$ .

Next, it follows from Lemma 3.6 that there exists a matrix  $\tilde{P}_i$  such that  $\|\tilde{P}_i\|_F = \mathcal{O}(1/(\sqrt{\eta}\alpha))$  and  $\|T'_4(\cdot, \cdot, \tilde{P}_i) - S'_i\|_F = \mathcal{O}(\sqrt{\eta/\alpha})$ . Furthermore, with probability at least 0.99, there exists a sequence of  $b_j \in [-D, D]$ , for  $j \in [4k]$ , such that  $T'_4(\cdot, \cdot, \sum_{j \in [4k]} b_j x^{(j)} \otimes y^{(j)}) = T'_4(\cdot, \cdot, \tilde{P}_i)$ .

Consider a cover,  $\mathcal{C}$ , of the interval  $[-D, D]$  with points spaced at intervals of length  $\tau = \mathcal{O}\left(\frac{\sqrt{\eta}}{\alpha k(\Delta^4 + 1/\alpha^4)}\right)$ . Since we uniformly sample  $a_j$ 's, with probability at least  $(\tau/D)^{\mathcal{O}(k)}$ , for all  $j \in [4k]$ ,  $|b_j - a_j| \leq \tau$ , and we condition on this event. Thus,

$$\begin{aligned} \left\| T'_4\left(\cdot, \cdot, \sum_{j \in [4k]} a_j x^{(j)} y^{(j)}\right) - S'_i \right\|_F &\leq \left\| T'_4\left(\cdot, \cdot, \sum_{j \in [4k]} b_j x^{(j)} \otimes y^{(j)}\right) - S'_i \right\|_F + \left\| T'_4\left(\cdot, \cdot, \sum_{j \in [4k]} (b_j - a_j) x^{(j)} \otimes y^{(j)}\right) \right\|_F \\ &\leq \mathcal{O}(\sqrt{\eta/\alpha}) + \mathcal{O}(\tau \Delta^2) \leq \mathcal{O}(\sqrt{\eta/\alpha}). \end{aligned} \tag{3.23}$$

Taking the linear combinations with coefficients  $a_j$  in Equation (3.22), we have

$$\begin{aligned} \hat{T}_4\left(\cdot, \cdot, \sum_{j \in [4k]} a_j x^{(j)} \otimes y^{(j)}\right) - S_i &= (\hat{T}_4 - T_4 + T'_4)\left(\cdot, \cdot, \sum_{j \in [4k]} a_j x^{(j)} \otimes y^{(j)}\right) - S'_i - \mu_i \otimes \mu_i \\ &\quad + \sum_{j \in [4k]} a_j \sum_{i \in [k]} w_i (S'_i x^{(j)}) \otimes (S'_i y^{(j)}) + \sum_{j \in [4k]} a_j \sum_{i \in [k]} w_i (S'_i y^{(j)}) \otimes (S'_i x^{(j)}) \\ &\quad + \sum_{j \in [4k]} a_j \sum_{i \in [k]} w_i (-2\mu_i^{\otimes 2} \langle \mu_i, x^{(j)} \rangle \langle \mu_i, y^{(j)} \rangle). \end{aligned} \tag{3.24}$$

Setting  $P_i = (\hat{T}_4 - T_4 + T'_4)\left(\cdot, \cdot, \sum_{j \in [4k]} a_j x^{(j)} y^{(j)}\right) - S'_i$ , it follows from Lemma 2.9 that with probability at least 0.99,  $(\hat{T}_4 - T_4)\left(\cdot, \cdot, \sum_{j \in [4k]} a_j x^{(j)} y^{(j)}\right)$  has Frobenius norm  $\mathcal{O}(kD\eta)$  and it follows from Equation (3.23) that with probability at least 0.99,  $T'_4\left(\cdot, \cdot, \sum_{j \in [4k]} a_j x^{(j)} y^{(j)}\right) - S'_i$  has Frobenius norm  $\mathcal{O}(\sqrt{\eta/\alpha})$ . Setting the remaining terms to  $Q_i$ , with probability at least 0.99 we can bound their

Frobenius norm as follows:

$$\begin{aligned}
\|Q_i\|_F &\leq \|\mu_i \otimes \mu_i\|_F + \left\| \left( \sum_{i \in [k]} w_i S'_i \oplus S'_i + w_i S'_i \ominus S'_i - 2w_i \mu_i^{\otimes 4} \right) \left( \cdot, \cdot, \sum_{j \in [4k]} a_j x^{(j)} y^{(j)} \right) \right\|_F \\
&\leq \frac{4}{\alpha} + \left( 2 \max_{i \in k} \|S'_i\|_F^2 + \frac{32}{\alpha^2} + k\tau \right) \cdot \|\tilde{P}\|_F \\
&\leq \frac{4}{\alpha} + \mathcal{O}\left(\frac{1}{\sqrt{\eta}\alpha} \left(\Delta + \frac{1}{\alpha}\right)^2\right) \\
&\leq \mathcal{O}\left(\frac{1 + \Delta^2}{\sqrt{\eta}\alpha^3}\right),
\end{aligned} \tag{3.25}$$

where the first inequality follows from the triangle inequality, the second follows from our assumptions that  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$ ,  $\sum_{j \in [4k]} b_j x^{(j)} y^{(j)} = \tilde{P}_i$  in the span of the  $S'_i$ , and  $|a_j - b_j| \leq \tau$  for all  $j \in [4k]$ , and the third inequality follows from the definition of  $S'_i$ , the bound on  $\|\tilde{P}\|_F$  and the bound on  $\|S_i - I\|_F$ .  $\square$

### 3.5 Finding a Low-dimensional Subspace for Exhaustive Search

In this subsection, we will prove Proposition 3.5.

We start by extending Theorem 4 of [MV10], which shows that large parameter distance between pairs of univariate Gaussian mixtures implies large distance between their low-degree moments. In the following, we use  $M_j(F) = \mathbb{E}[F^j]$  to denote the  $j$ -th moment of a distribution  $F$ . We show:

**Lemma 3.9.** *There exists a constant  $C > 0$  such that the following holds: Fix any  $D > 0$  and  $0 \leq \beta \leq 1/(2(2k-1)!D^{2k-3})$ . Suppose that  $F = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \sigma_i^2)$  is a univariate  $k$ -mixture of Gaussians with  $w_i \geq \beta$ , and  $|\mu_i|, \sigma_i \leq D$ , for all  $i \in [k]$ . If  $|\mu_i| + |\sigma_i^2 - 1| \geq \beta$  for some  $i \leq k$ , then*

$$\max_{j \in [2k]} |M_j(F) - M_j(\mathcal{N}(0, 1))| \geq \beta^{C^{k+1}(k+1)!-1}.$$

We give the proof of Lemma 3.9 in Section 3.6.

**Lemma 3.10** (Bounding  $\mu_i$ 's and  $S_i$ 's in non-influential directions for  $\mathbb{E}[h_m(\mathcal{M})]$ ). *Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians on  $\mathbb{R}^d$  satisfying  $w_i \geq \alpha$ ,  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$ , and  $\|\Sigma_i - I\|_F \leq \Delta$  for every  $i \in [k]$ . For some  $B \in \mathbb{R}$ , let  $u \in \mathbb{R}^d$  be a unit vector such that  $|\mathbb{E}[h_m(\langle \mathcal{M}, u \rangle)]| \leq B$  for all  $m \in [2k]$ . Then, for  $\delta = 2^{O(k)} B^{1/(C^{k+1}(k+1)!)}$  and  $S_i = \Sigma_i - I$ , we have that:*

1. for all  $i \leq k$ ,  $|\langle u, \mu_i \rangle|, |u^\top (I - \Sigma_i) u| \leq \delta$ ,
2.  $\|S_i u\|_2^2 \leq 20\delta\Delta/\alpha^2 + B/\alpha$ ,

where  $C > 0$  is a fixed universal constant.

*Proof.* The 1-D random variable  $\langle u, \mathcal{M} \rangle$  is a mixture of Gaussians described by  $\sum_{i=1}^k w_i \mathcal{N}(\langle \mu_i, u \rangle, u^\top \Sigma_i u)$ . Towards a contradiction, assume that there is an  $i \in [k]$  such that

$|\langle u, \mu_i \rangle| + |u^\top (I - \Sigma_i)u| \geq \delta$ . Then, applying Lemma 3.9, yields that there is a  $j \in [2k]$  such that  $|M_j(\langle u, \mathcal{M} \rangle) - M_j(\mathcal{N}(0, 1))| \geq \delta^{C^{k+1}(k+1)!-1}$ . Applying Fact 2.5 implies that there exists an  $m \in [2k]$  such that

$$\left| \mathbb{E} [h_m(\langle u, \mathcal{M} \rangle)] \right| > 2^{-O(k)} \delta^{C^{k+1}(k+1)!-1} \gg B ,$$

yielding a contradiction.

We can now prove the second part. Recall that for  $S_i = \Sigma_i - I$  for every  $i$ , we have that

$$\mathbb{E} [h_4(\mathcal{M})] = \sum_{i=1}^k w_i \text{Sym} (3(S_i \otimes S_i) + 6(S_i \otimes \mu_i^{\otimes 2}) + \mu_i^{\otimes 4}) .$$

We consider the  $d \times d$  matrix obtained by the natural flattening of the  $d \times d$  tensor  $u^{\otimes 2} \cdot \mathbb{E} [h_4(\mathcal{M})]$ . Then, we can write:

$$\begin{aligned} u^{\otimes 2} \cdot \mathbb{E} [h_4(\mathcal{M})] &= \sum_{i=1}^k w_i \left( (u^\top S_i u) S_i + 2(S_i u)(S_i u)^\top + \langle u, \mu_i \rangle^2 S_i \right. \\ &\quad \left. + 2\langle u, \mu_i \rangle \mu_i (S_i u)^\top + 2\langle u, \mu_i \rangle (S_i u) \mu_i^\top + (u^\top S_i u) \mu_i \mu_i^\top + \langle u, \mu_i \rangle^2 \mu_i \mu_i^\top \right) . \end{aligned} \quad (3.26)$$

Now, from the first part, we know that for all  $i \in [k]$ ,  $|u^\top S_i u| \leq \delta$  and the hypothesis of the lemma gives us that  $\|S_i\|_F = \|\Sigma_i - I\|_F \leq \Delta$ . Thus, for each  $i$ , the first term in the summation above has Frobenius norm at most  $\Delta\delta$ . Using that  $\langle u, \mu_i \rangle^2 \leq \delta^2$  from the first part of the lemma, yields that, for each  $i$ , the Frobenius norm of the third term is at most  $\Delta\delta^2$ .

Next, using in addition that  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$  yields that, for each  $i$ , the Frobenius norm of the 4th and 5th terms are at most  $2\delta\Delta/\sqrt{\alpha}$  and the Frobenius norm of the 6th and 7th terms are at most  $\delta/\alpha$ . Thus, for each  $i$  and all but the 2nd term in the summation above, we have an upper bound on the Frobenius norm of  $4\delta\Delta/\alpha$ .

Now, since  $|\mathbb{E} [h_4(\langle \mathcal{M}, u \rangle)]| \leq B$ , and  $u$  is a unit vector, we have that  $\|u^{\otimes 2} \mathbb{E} [h_4(\mathcal{M})]\|_F \leq B$ . Thus, combining the aforementioned argument with the triangle inequality, we have for each  $i$ ,

$$\begin{aligned} \|S_i u\|_2^2 &= \|S_i u (S_i u)^\top\|_F \leq \frac{1}{\alpha} \left\| u^{\otimes 2} \cdot \mathbb{E} [h_4(\mathcal{M})] \right\|_F + \sum_{i \in [k]} w_i \left( (u^\top S_i u + \langle u, \mu_i \rangle^2) \|S_i\|_F \right) \\ &\quad + \sum_{i \in [k]} 4w_i \left( \langle u, \mu_i \rangle \|\mu_i\|_2 \|S_i u\|_2 \right) + \sum_{i \in [k]} 4w_i \left( \langle u, \mu_i \rangle^2 + u^\top S_i u \right) \|\mu_i\|_2^2 \\ &\leq B/\alpha + 15\delta\Delta/\alpha , \end{aligned}$$

and the claim follows.  $\square$

**Lemma 3.11** (Subspace covering all the means and large singular vectors of  $S_i = \Sigma_i - I$ ). *Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians on  $\mathbb{R}^d$  satisfying  $w_i \geq \alpha$ ,  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$ , and  $\|\Sigma_i - I\|_F \leq \Delta$  for all  $i \in [k]$ . Given  $0 < \eta < 1$ , let  $\hat{T}_m$  satisfy  $\|\hat{T}_m - \mathbb{E} [h_m(\mathcal{M})]\|_F \leq \eta$  for every  $m \in [4k]$  and let  $\lambda \geq 2\eta$ . Let  $V$  be the span of all the left singular vectors of the  $d \times d^{m-1}$  matrix obtained by the natural flattening of  $\hat{T}_m$  with singular values at least  $\lambda$ . Then, for  $\delta = 2\lambda^{1/(2C^{k+1}(k+1)!)}$ , we have that:*

$$1. \dim(V) \leq (4k\eta^2 + k^{O(k)}) \mathcal{O}(1 + 1/\alpha + \Delta)^{4k} / \lambda^2,$$

2. Let

$$V_{\text{inf}} = \{\mu_i\}_{i \in [k]} \cup \left\{ v \mid \exists i \in [k], \text{ s.t. } \|v\|_2 = 1 \text{ and } v \text{ is an eigenvector of } S_i \text{ and } \|S_i v\|_2 \geq \sqrt{\delta} \right\}_{i \leq k}.$$

Then, for every unit vector  $v \in V_{\text{inf}}$ ,  $\|v - \Pi_V v\|_2^2 \leq 20\delta^{1/4}\Delta/\alpha^2$ , where  $\Pi_V v$  is the projection of  $v$  onto  $V$ .

*Proof.* From Fact 2.6, we have that  $\mathbb{E}[h_m(\mathcal{M})] = \sum_{i \in [k]} w_i \mathbb{E}[h_m(G_i)]$ , where  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$ , and since  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$  and  $\|\Sigma_i - I\|_F \leq \Delta$ , it follows that  $\|\mathbb{E}[h_m(\mathcal{M})]\|_F^2 \leq \mathcal{O}(m(1 + 1/\alpha + \Delta))^{4m}$ . From Proposition 3.3, we know that

$$\left\| \hat{T}_m \right\|_F^2 \leq 2 \left\| \hat{T}_m - \mathbb{E}[h_m(\mathcal{M})] \right\|_F^2 + 2 \left\| \mathbb{E}[h_m(\mathcal{M})] \right\|_F^2 \leq \eta^2 + \mathcal{O}(m(1 + 1/\alpha + \Delta))^{4m}.$$

Thus, the number of singular vectors of the  $d \times d^{m-1}$  flattening of  $\hat{T}_m$  with a singular value  $\geq \lambda$  is at most  $(\eta^2 + \mathcal{O}(m(1 + 1/\alpha + \Delta))^{4m})/\lambda^2$ . Summing up this bound for all  $m \in [4k]$ , yields the claimed upper bound on  $\dim(V)$ .

For the second part, we will first bound  $\langle u, v \rangle$  for any unit vector  $u$  orthogonal to the subspace  $V$ . Towards this, observe that since  $u$  is orthogonal to  $V$  and  $\|u\|_2 = 1$ , we have

$$\left\| u \cdot \mathbb{E}[h_m(\mathcal{M})] \right\|_F \leq \left\| u \hat{T}_m \right\|_F + \left\| \hat{T}_m - \mathbb{E}[h_m(\mathcal{M})] \right\|_F \leq \lambda + \eta \leq 2\lambda,$$

where  $u \cdot \mathbb{E}[h_m(\mathcal{M})]$  is a matrix-vector product of  $u$  with a  $d \times d^{m-1}$  flattening of  $\mathbb{E}[h_m(\mathcal{M})]$ . For  $\delta = 2\lambda^{1/(C^{k+1}(k+1)!)}$ , applying Lemma 3.10 yields that

$$\langle \mu_i, u \rangle^2 + \|S_i u\|_2^2 \leq \delta^2 + 20\delta\Delta/\alpha \leq 20\delta\Delta/\alpha^2. \quad (3.27)$$

Now, if  $v$  is one of the  $\mu_i$ 's, then we immediately get from Equation 3.27 that  $\langle v, u \rangle^2 \leq 20\delta\Delta/\alpha^2$ . Similarly, note that if  $v$  is a unit length eigenvector of  $S_i$  satisfying  $\|S_i v\|_2^2 \geq \sqrt{\delta}$ , then,

$$\langle u, v \rangle^2 = \frac{1}{\|S_i v\|_2^2} \langle u, S_i v \rangle^2 = \frac{1}{\|S_i v\|_2^2} \langle S_i u, v \rangle^2 \leq \frac{\|S_i u\|_2^2}{\|S_i v\|_2^2}.$$

In both cases, setting  $u = (v - \Pi_V v)/\|v - \Pi_V v\|_2$  completes the proof.  $\square$

We can now complete the proof of Proposition 3.5:

*Proof of Proposition 3.5.* We know that  $\hat{S}_i - P_i - S_i$  is a symmetric, rank- $k'$  matrix such that  $k' = \mathcal{O}(k^2)$ , described by the eigenvalue decomposition  $\sum_{i=1}^{k'} \tau_i v_i v_i^\top$ , where  $v_i$ 's are the eigenvectors and  $\tau_i$ 's are the corresponding eigenvalues. Since  $\|S_i\|_F \leq \Delta$  and

$$\left\| \hat{S}_i \right\|_F \leq \|P_i\|_F + \|Q_i\|_F + \|S_i\|_F \leq \mathcal{O}\left(\sqrt{\eta/\alpha}\right) + \mathcal{O}\left(\frac{1 + \Delta^2}{\sqrt{\eta}\alpha^3}\right) + \Delta = \mathcal{O}\left(\frac{1 + \Delta^2}{\sqrt{\eta}\alpha^3}\right),$$

we have that the number of singular values of  $\hat{S}_i$  that exceed  $\delta^{1/4}$  is at most  $\mathcal{O}\left(\frac{1+\Delta^2}{\sqrt{\eta}\alpha^3\sqrt{\delta}}\right)$ . Recall that from Lemma 3.11 it follows that the dimension of the subspace  $V$  is at most  $k^{\mathcal{O}(k)}\mathcal{O}(1+1/\alpha+\Delta)^{4k}/\lambda^2$ . Thus, the dimension of  $V'$  is at most

$$k^{\mathcal{O}(k)}\mathcal{O}\left(\frac{(1+1/\alpha+\Delta)^{4k}}{\lambda^2}\right) + \mathcal{O}\left(\frac{1+\Delta^2}{\sqrt{\eta}\alpha^3\sqrt{\delta}}\right) = \mathcal{O}\left(\frac{k^{\mathcal{O}(k)}(1+\frac{1}{\alpha}+\Delta)^{4k+5}}{\eta^2}\right).$$

Since  $V'$  contains  $V$  constructed in Lemma 3.11, we immediately obtain that for every  $\mu_i$ ,  $\|\mu_i - \Pi_{V'}\mu_i\|_2^2 \leq \frac{20}{\alpha^2}\sqrt{\delta}\Delta$ .

Next, let  $u$  be a unit vector orthogonal to  $V'$ . Then, since  $V'$  contains the  $V$  described in Lemma 3.11, we know that  $\|S_i u\|_2^2 \leq \frac{20}{\alpha^2}\sqrt{\delta}\Delta$ . Similarly, since  $V'$  contains all eigenvectors of  $\hat{S}_i$  with singular values exceeding  $\delta^{1/4}$ , we know that  $\|\hat{S}_i u\|_2^2 \leq \delta^{1/2}$ . Thus, we can conclude that  $\|(\hat{S}_i - S_i)u\|_2^2 \leq \frac{100}{\alpha^2}\sqrt{\delta}\Delta$ . Let  $Q_i = \sum_{j=1}^{k'} \tau_j v_j v_j^\top$  with orthonormal  $v_j \in \mathbb{R}^d$ . We know such  $\tau_j$ 's and  $v_j$ 's exist because of the upper bound on  $\text{rank}(Q_i)$ . Therefore, for any  $j$ ,  $|v_j^\top(\hat{S}_i - S_i)u| \leq \frac{10}{\alpha}\delta^{1/4}\Delta^{1/2}$ . On the other hand, for any  $j$ , we have that

$$v_j^\top(\hat{S}_i - S_i)u \geq \langle v_j, u \rangle \tau_j - \|P_i\|_F = \langle v_j, u \rangle \tau_j - \mathcal{O}(\sqrt{\eta}).$$

Combining the two bounds above, yields that whenever  $\tau_j \geq \delta^{1/4}$ ,

$$|\langle v_j, u \rangle| \leq \mathcal{O}(\sqrt{\eta}/\tau_j) + \frac{10}{\alpha\tau_j}\delta^{1/4}\Delta^{1/2} \leq \frac{10}{\alpha}\delta^{1/2}\Delta^{1/2}.$$

Thus, the matrix  $\hat{Q}_i = \sum_{j=1}^{k'} \tau_j \Pi_{V'} v_j (\Pi_{V'} v_j)^\top$  has its range space in  $V'$  and satisfies

$$\|\hat{Q}_i - Q_i\|_F \leq \mathcal{O}(k^2\delta^{1/4}) + \mathcal{O}\left(\frac{k^2}{\alpha}\delta^{1/2}\Delta^{1/2}\right) = \mathcal{O}\left(\frac{k^2}{\alpha}\delta^{1/2}\Delta^{1/2}\right).$$

□

### 3.6 Parameter vs Moment Distance for Gaussian Mixtures

In this subsection, we prove Lemma 3.9. To that end, we will use the following two results; the second one is from [MV10].

**Lemma 3.12.** *Suppose  $\mathcal{N}(\mu_1, \sigma_1^2)$  and  $\mathcal{N}(\mu_2, \sigma_2^2)$  are univariate Gaussians with  $|\mu_i|, |\sigma_i| \leq D$ , for some  $D \in \mathbb{R}_+$ . If  $|\mu_1 - \mu_2| + |\sigma_1^2 - \sigma_2^2| \leq \beta$ , then the distance between raw moments of two Gaussians is*

$$|M_j(\mathcal{N}(\mu_1, \sigma_1^2)) - M_j(\mathcal{N}(\mu_2, \sigma_2^2))| \leq (j+1)!D^{j-1}\beta.$$

*Proof.* By Proposition 2.3, the  $j$ -th raw moment of a Gaussian  $\mathcal{N}(\mu, \sigma^2)$  is a sum of monomials in  $\mu$  and  $\sigma^2$  of degree  $j$ . There are at most  $(j+1)!$  terms in the polynomial. Thus, changing the mean or the variance by at most  $\beta$  will change the  $j$ -th moment by at most  $(j+1)!D^{j-1}\beta$ . □

**Theorem 3.13.** ([MV10]) Let  $F, F'$  be two univariate mixtures of Gaussians:  $F = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \sigma_i^2)$  and  $F' = \sum_{i=1}^{k'} w'_i \mathcal{N}(\mu'_i, \sigma'^2_i)$ . There is a constant  $c > 0$  such that, for any  $\beta < c$ , if  $F, F'$  satisfy:

1.  $w_i, w'_i \in [\beta, 1]$
2.  $|\mu_i|, |\mu'_i| \leq 1/\beta$
3.  $|\mu_i - \mu_{i'}| + |\sigma_i^2 - \sigma_{i'}^2| \geq \beta$  and  $|\mu'_i - \mu'_{i'}| + |\sigma'^2_i - \sigma'^2_{i'}| \geq \beta$  for all  $i \neq i'$
4.  $\beta \leq \min_{\pi} \sum_i (|w_i - w'_{\pi(i)}| + |\mu_i - \mu'_{\pi(i)}| + |\sigma_i^2 - \sigma'^2_{\pi(i)}|)$ , where the minimization is taken over all mappings  $\pi : \{1, \dots, k\} \rightarrow \{1, \dots, k'\}$ ,

then

$$\max_{j \in [2(k+k'-1)]} |M_j(F) - M_j(F')| \geq \beta^{O(k)}.$$

We are now ready to complete the proof of Lemma 3.9.

*Proof of Lemma 3.9.* We proceed via induction on  $k$ . Consider the base case, i.e.,  $k = 1$ . Then, either  $|\mu_1| \geq \beta/2$  or  $|\sigma_1 - 1| \geq \beta/2$ , and thus the first or second moment differ by at least  $\beta^2/4$ . Let the inductive hypothesis be that Lemma 3.9 holds for at most  $k$  components.

Consider the case where  $|\mu_i - \mu_{i'}| + |\sigma_i^2 - \sigma_{i'}^2| \geq \beta^{C^k k!}$  for all pairs of components  $i, i' \in [k]$ . Then, by Theorem 3.13, we have that

$$\max_{j \in [2k]} |M_j(F) - M_j(\mathcal{N}(0, 1))| \geq \beta^{C^{k+1} k!} \geq \beta^{C^{k+1} (k+1)!-1},$$

and the lemma follows.

Otherwise, we know that there exists a pair of components with parameter distance less than  $\beta^{C^k k!}$ . In this case, we merge these two components and get a  $(k-1)$ -mixture  $F'$ . By Lemma 3.12, the distance between the  $j$ -th moments of  $F'$  and  $F$  is at most  $(j+1)! D^{j-1} \beta^{C^k k!}$ . Since we still have  $|\mu'_i| + |\sigma'^2_i - 1| \geq \beta - 3\beta^{C^k k!}$  for all components  $i$  in  $F'$ , the inductive hypothesis implies that

$$\max_{j \in [2k-2]} |M_j(F') - M_j(\mathcal{N}(0, 1))| \geq \left( \beta - 3\beta^{C^k k!} \right)^{C^k (k-1)!-1}.$$

By the triangle inequality, we can write

$$\begin{aligned} \max_{j \in [2k]} |M_j(F) - M_j(\mathcal{N}(0, 1))| &\geq \max_{j \in [2k-2]} |M_j(F') - M_j(\mathcal{N}(0, 1))| - \max_{j \in [2k-2]} |M_j(F) - M_j(F')| \\ &\geq \left( \beta - 3\beta^{C^k k!} \right)^{C^k (k-1)!-1} - (2k-1)! D^{2k-3} \beta^{C^k k!} \\ &\geq \beta^{C^{k+1} (k+1)!-1}. \end{aligned}$$

The last inequality follows from the assumption that  $\beta \leq 1/(2(2k-1)! D^{2k-3})$ . This completes the proof of Lemma 3.9.  $\square$

## 4 Robust Partial Cluster Recovery

In this section, we give two robust *partial clustering* algorithms. A partial clustering algorithm takes a set of points  $X = \cup_{i \leq k} X_i$  with true clusters  $X_1, X_2, \dots, X_k$  and outputs a partition of the sample  $X = X'_1 \cup X'_2$  such that  $X'_1 = \cup_{i \in S} X_i$  and  $X'_2 = \cup_{i \notin S} X_i$ , for some subset  $S \subseteq [k]$  of size  $1 \leq |S| < k$ . That is, a partial clustering algorithm partitions the sample into two non-empty parts so that each part is a sample from a “sub-mixture”. This is a weaker guarantee than clustering the entire mixture, which must find each of the original  $X_i$ ’s. We show that the relaxed guarantee is feasible even when the mixture as a whole is not clusterable. In our setting, we will get an approximate (that is, a small fraction of points are misclassified) partial clustering that works for  $\varepsilon$ -corruptions  $Y$  of any i.i.d. sample  $X$  from a mixture of  $k$  Gaussians, as long as there is a pair of components in the original mixture that have large total variation distance between them.

A partial clustering algorithm such as above was one of the innovations in [BK20b] that allowed for a polynomial-time algorithm for clustering all fully clusterable Gaussian mixtures.

In this section, we build on the ideas in [BK20b] to derive two new partial clustering algorithms that work even when the original mixture is not fully clusterable. Both upgrade the results of [BK20b] by handling mixtures with arbitrary weights  $w_i$ ’s instead of uniform weights and handling mixtures where not all pairs of components are well-separated in TV distance. The first algorithm succeeds under the information-theoretically minimal separation assumption (i.e. separation in total variation distance) but runs in time exponential in the inverse mixing weight. The second algorithm is a key innovation of this paper – it gives an algorithm that runs in polynomial time in the inverse mixing weight at the cost of handling separation only in relative Frobenius distance. This improved running time guarantee (at the cost of strong separation requirement that we mitigate through a novel standalone spectral separation step in Section 5) is crucial to obtaining the fully polynomial running time in our algorithm.

In order to state the guarantees of our algorithms, we first formulate a notion of parameter separation (same as the one employed in [BK20b, DHKK20]) as the next definition.

**Definition 4.1** ( $\Delta$ -Parameter Separation). We say that two Gaussian distributions  $\mathcal{N}(\mu_1, \Sigma_1)$  and  $\mathcal{N}(\mu_2, \Sigma_2)$  are  $\Delta$ -parameter separated if at least one of the following three conditions hold:

1. **Mean-Separation:**  $\exists v \in \mathbb{R}^d$  such that  $\langle \mu_1 - \mu_2, v \rangle^2 > \Delta^2 v^\top (\Sigma_1 + \Sigma_2) v$ ,
2. **Spectral-Separation:**  $\exists v \in \mathbb{R}^d$  such that  $v^\top \Sigma_1 v > \Delta v^\top \Sigma_2 v$ ,
3. **Relative-Frobenius Separation:**  $\Sigma_i$  and  $\Sigma_j$  have the same range space and  $\left\| \Sigma_1^{\dagger/2} (\Sigma_2 - \Sigma_1) \Sigma_1^{\dagger/2} \right\|_F^2 > \Delta^2 \left\| \Sigma_1^\dagger \Sigma_2 \right\|_{\text{op}}^2$ .

As shown in [DHKK20, BK20b], if a pair of Gaussians is  $(1 - \exp(-O(\Delta \log \Delta)))$ -separated in total variation distance, then, they are  $\Delta$ -parameter separated.

Our first algorithm succeeds in robust partial clustering whenever there is a pair of component Gaussians that are  $\Delta$ -parameter separated. The running time of this algorithm grows exponentially in the reciprocal of the minimum weight in the mixture.



**Theorem 4.2** (Robust Partial Clustering in TV Distance). *Let  $0 \leq \varepsilon < \alpha \leq 1$ , and  $\eta > 0$ . There is an algorithm with the following guarantees: Let  $\{\mu_i, \Sigma_i\}_{i \leq k}$  be means and covariances of  $k$  unknown Gaussians. Let  $Y$  be an  $\varepsilon$ -corruption of a sample  $X$  of size  $n \geq (dk)^{Ct} / \varepsilon$  for a large enough constant  $C > 0$ , from  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  satisfying Condition 2.47 with parameters  $t = (k/\eta)^{O(k)}$  and  $\gamma \leq \varepsilon d^{-8t} k^{-Ct}$ , for a sufficiently large constant  $C > 0$ . Suppose further that  $w_i \geq \alpha > 2\varepsilon$  for every  $i$  and that there are  $i, j$  such that  $\mathcal{N}(\mu_i, \Sigma_i)$  and  $\mathcal{N}(\mu_j, \Sigma_j)$  are  $\Delta$ -parameter separated for  $\Delta = (k/\eta)^{O(k)}$ .*

*Then, the algorithm on input  $Y$ , runs in time  $n^{(k/\eta)^{O(k)}}$ , and with probability at least  $2^{-O\left(\frac{1}{\alpha} \log\left(\frac{k}{\eta\alpha}\right)\right)}$  over the draw of  $X$  and the algorithm's random choices, the algorithm outputs a partition of  $Y$  into  $Y_1, Y_2$  satisfying:*

1. **Partition respects clustering:** for each  $i$ ,  $\max\left\{\frac{k}{n}|Y_1 \cap X_i|, \frac{k}{n}|Y_2 \cap X_i|\right\} \geq 1 - \eta - O(\varepsilon/\alpha^4)$ , and,
2. **Partition is non-trivial:**  $\max_i \frac{k}{n}|X_i \cap Y_1|, \max_i \frac{k}{n}|X_i \cap Y_2| \geq 1 - \eta - O(\varepsilon/\alpha^4)$ .

Our proof of the above theorem is based on a relatively straightforward extension of the ideas of [BK20b], albeit with two key upgrades 1) allowing the input mixtures to have arbitrary mixing weights (at an exponential cost in the inverse of the minimum weight) and 2) handling mixtures where some pair of components may not be well-separated in TV distance.

In order to get our main result that gives a fully polynomial algorithm (including in the inverse mixing weights), we will use an incomparable variant of the above partial clustering method that only handles a weaker notion of parameter separation, but runs in fixed polynomial time.

**Theorem 4.3** (Robust Partial Clustering in Relative Frobenius Distance). *Let  $0 \leq \varepsilon < \alpha/k \leq 1$  and  $t \in \mathbb{N}$ . There is an algorithm with the following guarantees: Let  $\{\mu_i, \Sigma_i\}_{i \leq k}$  be means and covariances of  $k$  unknown Gaussians. Let  $Y$  be an  $\varepsilon$ -corruption of a sample  $X$  of size  $n \geq (dk)^{Ct} / \varepsilon$  for a large enough constant  $C > 0$ , from  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  that satisfies Condition 2.47 with parameters  $2t$  and  $\gamma \leq \varepsilon d^{-8t} k^{-Ck}$ , for a large enough constant  $C > 0$ . Suppose further that  $w_i \geq \alpha > 2\varepsilon$  for each  $i \in [k]$ , and that for some  $t \in \mathbb{N}$ ,  $\beta > 0$  there exist  $i, j \leq k$  such that  $\|\Sigma^{+1/2}(\Sigma_i - \Sigma_j)\Sigma^{+1/2}\|_F^2 = \Omega((k^2 t^4)/(\beta^{2/t} \alpha^4))$ , where  $\Sigma$  is the covariance of the mixture  $\mathcal{M}$ . Then, the algorithm runs in time  $n^{O(t)}$ , and with probability at least  $2^{-O\left(\frac{1}{\alpha} \log\left(\frac{k}{\eta}\right)\right)}$  over the random choices of the algorithm, outputs a partition  $Y = Y_1 \cup Y_2$  satisfying:*

1. **Partition respects clustering:** for each  $i$ ,  $\max\left\{\frac{1}{w_i n}|Y_1 \cap X_i|, \frac{1}{w_i n}|Y_2 \cap X_i|\right\} \geq 1 - \beta - O(\varepsilon/\alpha^4)$ , and,
2. **Partition is non-trivial:**  $\max_i \frac{1}{w_i n}|X_i \cap Y_1|, \max_i \frac{1}{w_i n}|X_i \cap Y_2| \geq 1 - \beta - O(\varepsilon/\alpha^4)$ .

The starting point for the proof of the above theorem is the observation that the running time of our first algorithm above is exponential in the inverse mixing weight almost entirely because of dealing with spectral separation (which requires the use of ‘‘certifiable anti-concentration’’ that we define in the next subsection). We formulate a variant of relative Frobenius separation (that is directly useful to us) and prove that whenever the original mixture has a pair of components separated in this notion, we can in fact obtain a fully polynomial partial clustering algorithm building on the ideas in [BK20b].

## 4.1 Algorithm

Our algorithm will solve SoS relaxations of a polynomial inequality system. The constraints here use the input  $Y$  to encode finding a sample  $X'$  (the intended setting being  $X' = X$ , the original uncorrupted sample) and a cluster  $\hat{C}$  in  $X'$  of size  $= \alpha n$ , indicated by  $z_i$ s (the intended setting is simply the indicator for any of the  $k$  true clusters) satisfying properties of Gaussian distribution (certifiable hypercontractivity and anti-concentration).

Covariance constraints introduce a matrix valued indeterminate  $\Pi$  intended to be the square root of  $\hat{\Sigma}$ , the sos variable for the covariance of a single component.

$$\text{Covariance Constraints: } \mathcal{A}_1 = \left\{ \begin{array}{l} \Pi = UU^\top \\ \Pi^2 = \hat{\Sigma} \end{array} \right\} \quad (4.1)$$

The intersection constraints force that  $X'$  be  $\varepsilon$ -close to  $Y$  (and thus,  $2\varepsilon$ -close to unknown sample  $X$ ).

$$\text{Intersection Constraints: } \mathcal{A}_2 = \left\{ \begin{array}{l} \forall i \in [n], \quad m_i^2 = m_i \\ \sum_{i \in [n]} m_i = (1 - \varepsilon)n \\ \forall i \in [n], \quad m_i(y_i - x'_i) = 0 \end{array} \right\} \quad (4.2)$$

The subset constraints introduce  $z$ , which indicates the subset  $\hat{C}$  intended to be the true clusters of  $X'$ .

$$\text{Subset Constraints: } \mathcal{A}_3 = \left\{ \begin{array}{l} \forall i \in [n]. \quad z_i^2 = z_i \\ \sum_{i \in [n]} z_i = \alpha n \end{array} \right\} \quad (4.3)$$

Parameter constraints create indeterminates to stand for the covariance  $\hat{\Sigma}$  and mean  $\hat{\mu}$  of  $\hat{C}$  (indicated by  $z$ ).

$$\text{Parameter Constraints: } \mathcal{A}_4 = \left\{ \begin{array}{l} \frac{1}{\alpha n} \sum_{i=1}^n z_i (x'_i - \hat{\mu}) (x'_i - \hat{\mu})^\top = \hat{\Sigma} \\ \frac{1}{\alpha n} \sum_{i=1}^n z_i x'_i = \hat{\mu} \end{array} \right\} \quad (4.4)$$

Certifiable Hypercontractivity :  $\mathcal{A}_4 =$

$$\left\{ \begin{array}{l} \forall t \leq 2s \quad \frac{1}{\alpha^2 n^2} \sum_{i,j \leq n} z_i z_j (Q(x'_i - x'_j) - \mathbb{E}_z Q)^{2t} \leq \left( \frac{1}{\alpha^2 n^2} \sum_{i,j \leq n} z_i z_j (Q(x'_i - x'_j) - \mathbb{E}_z Q)^2 \right)^t \\ \frac{1}{\alpha^2 n^2} \sum_{i,j \leq n} z_i z_j (Q(x'_i - x'_j) - \mathbb{E}_z Q)^2 \leq \frac{6}{\alpha^2} \|Q\|_F^2 \end{array} \right\} \quad (4.5)$$

Here, we used the shorthand  $\mathbb{E}_z Q = \frac{1}{\alpha^2 n^2} \sum_{i,j \leq n} z_i z_j Q(x'_i - x'_j)$ .

In the constraint system for our first algorithm, we will use the following certifiable anti-concentration constraints on  $\hat{C}$  for  $\delta = \alpha^{-\text{poly}(k)}$  and  $\tau = \alpha/\text{poly}(k)$  and  $s(u) = 1/u^2$  for every

u.

$$\text{Certifiable Anti-Concentration : } \mathcal{A}_5 = \left\{ \begin{array}{l} \frac{1}{\alpha^2 n^2} \sum_{i,j=1}^n z_i z_j q_{\delta, \Sigma}^2 \left( \left( x'_i - x'_j \right), v \right) \leq 2^{s(\delta)} C \delta \left( v^\top \Sigma v \right)^{s(\delta)} \\ \frac{1}{\alpha^2 n^2} \sum_{i,j=1}^n z_i z_j q_{\tau, \Sigma}^2 \left( \left( x'_i - x'_j \right), v \right) \leq 2^{s(\tau)} C \tau \left( v^\top \Sigma v \right)^{s(\tau)} \end{array} \right\} \quad (4.6)$$

We note that the constraint system for our second algorithm (running in fixed polynomial time), we will not use  $\mathcal{A}_5$ . Towards proving Theorems 4.3 and 4.2 we use the following algorithm that differs only in the degree of the pseudo-distribution computed and the constraint system that the pseudo-distribution satisfies.

**Algorithm 4.4** (Partial Clustering).

**Given:** A sample  $Y$  of size  $n$ . An outlier parameter  $\varepsilon > 0$  and an accuracy parameter  $\eta > 0$ .

**Output:** A partition of  $Y$  into partial clustering  $Y_1 \cup Y_2$ .

**Operation:**

1. **SDP Solving:** Find a pseudo-distribution  $\tilde{\zeta}$  satisfying  $\cup_{i=1}^5 \mathcal{A}_i$  ( $\cup_{i=1}^4 \mathcal{A}_i$  for Theorem 4.3) such that  $\tilde{\mathbb{E}}_{\tilde{\zeta}} z_i \leq \alpha + o_d(1)$  for every  $i$ . If no such pseudo-distribution exists, output fail.
2. **Rounding:** Let  $M = \tilde{\mathbb{E}}_{z \sim \tilde{\zeta}} [z z^\top]$ .
  - (a) Choose  $\ell = O\left(\frac{1}{\alpha} \log(k/\eta)\right)$  rows of  $M$  uniformly at random and independently.
  - (b) For each  $i \leq \ell$ , let  $\hat{C}_i$  be the indices of the columns  $j$  such that  $M(i, j) \geq \eta^2 \alpha^5 / k$ .
  - (c) Choose a uniformly random  $S \subseteq [\ell]$  and output  $Y_1 = \cup_{i \in S} \hat{C}_i$  and  $Y_2 = Y \setminus Y_1$ .

## 4.2 Analysis

**Simultaneous Intersection Bounds.** The key observation for proving the first theorem is the following lemma that gives a sum-of-squares proof that no  $z$  that satisfies the constraints  $\cup_{i=1}^5 \mathcal{A}_i$  can have simultaneously large intersections with the  $\Delta$ -parameter separated component Gaussians.

**Lemma 4.5** (Simultaneous Intersection Bounds for TV-separated case). *Let  $Y$  be an  $\varepsilon$ -corruption of a sample  $X$  of size  $n \geq (dk)^{Ct} / \varepsilon$  for a large enough constant  $C > 0$ , from  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  satisfying Condition 2.47 with parameters  $t = (k/\eta)^{O(k)}$  and  $\gamma \leq \varepsilon d^{-8t} k^{-Ct}$ , for a sufficiently large constant  $C > 0$ . Suppose further that  $w_i \geq \alpha > 2\varepsilon$  for every  $i$  and that there are  $i, j$  such that  $\mathcal{N}(\mu_i, \Sigma_i)$  and  $\mathcal{N}(\mu_j, \Sigma_j)$  are  $\Delta$ -parameter separated for  $\Delta = (k/\eta)^{O(k)}$ . Then, there exists a partition of  $[k]$  into  $S \cup L$  such that,  $|S|, |L| < k$  and for  $z(X_r) = \frac{1}{w_r n} \sum_{i \in X_r} z_i$ ,*

$$\left\{ \cup_{i=1}^5 \mathcal{A}_i \mid \frac{z}{(k/\eta\alpha)^{\text{poly}(k)}} \left\{ \sum_{i \in S, j \in L} z(X_i) z(X_j) \leq O(k^2 \varepsilon / \alpha) + \eta / \alpha \right\} \right\}.$$

The proof of Lemma 4.5 is given in Section 4.3.

For the second theorem, we use the following version that strengthens the separation assumption and lowers the degree of the sum-of-squares proof (and consequently the running time of the algorithm) as a result.

**Lemma 4.6** (Simultaneous Intersection Bounds for Frobenius Separated Case). *Let  $X$  be a sample of size  $n \geq (dk)^{Ct}/\varepsilon$  for a large enough constant  $C > 0$ , from  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  that satisfies Condition 2.47 with parameters  $2t$  and  $\gamma \leq \varepsilon d^{-8t} k^{-Ck}$ , for a large enough constant  $C > 0$ . Suppose further that  $w_i \geq \alpha > 2\varepsilon$  for each  $i \in [k]$ , and that for some  $t \in \mathbb{N}$ ,  $\beta > 0$  there exist  $i, j \leq k$  such that  $\|\Sigma^{t/2}(\Sigma_i - \Sigma_j)\Sigma^{t/2}\|_F^2 = \Omega((k^2 t^4)/(\beta^{2/t} \alpha^2))$ , where  $\Sigma$  is the covariance of the mixture  $\mathcal{M}$ . Then, for any  $\varepsilon$ -corruption  $Y$  of  $X$ , there exists a partition of  $[k] = S \cup T$  such that*

$$\left\{ \bigcup_{i=1}^4 \mathcal{A}_i \right\} \Big|_{\frac{z}{2t}} \left\{ \sum_{i \in S} \sum_{j \in T} z(X_i)z(X_j) \leq O(k^2)\beta + O(k^2)\varepsilon/\alpha \right\}.$$

Here,  $z(X_r) = \frac{1}{w_r n} \sum_{i \in X_r} z_i$  for every  $r$ .

The proof of Lemma 4.6 is given in Section 4.4.

Notice that the main difference between the above two lemmas is the constraint systems they use. Specifically, the second lemma does *not* enforce certifiable anti-concentration constraints. As a result, there is a difference in the degree of the sum-of-squares proofs they claim; the degree of the SoS proof in the second lemma does not depend on the inverse minimum mixture weight.

First, we complete the proof of the Theorem 4.2. The proof of Theorem 4.3 is exactly the same except for the use of Lemma 4.6 (and thus has the exponent in the running time independent of  $1/\alpha$ ) instead of Lemma 4.5.

*Proof of Theorem 4.2.* Let  $\eta' = O(\eta^2 \alpha^3/k)$ . We will prove that whenever  $\Delta \geq \text{poly}(k/\eta')^k = \text{poly}\left(\frac{k}{\eta\alpha}\right)^k$ , Algorithm 4.4, when run with input  $Y$ , with probability at least 0.99, recovers a collection  $\hat{C}_1, \hat{C}_2, \dots, \hat{C}_\ell$  of  $\ell = O(\frac{1}{\alpha} \log k/\eta)$  subsets of indices satisfying  $|\cup_{i \leq \ell} \hat{C}_i| \geq (1 - \eta'/k^{40})n$  such that there is a partition  $S \cup L = [\ell]$ ,  $0 < |S| < \ell$  satisfying:

$$\min \left\{ \frac{1}{\alpha n} |\hat{C}_i \cap \cup_{j \in S} X_j|, \frac{1}{\alpha n} |\hat{C}_i \cap \cup_{j \in L} X_j| \right\} \leq 100\eta'/\alpha^3 + O(\varepsilon/\alpha^4). \quad (4.7)$$

We first argue that this suffices to complete the proof. Split  $[\ell]$  into two groups  $G_S, G_L$  as follows. For each  $i$ , let  $j = \text{argmax}_{r \in [\ell]} \frac{1}{\alpha n} |\hat{C}_i \cap X_r|$ . If  $j \in S$ , add it to  $G_S$ , else add it to  $G_L$ . Observe that this process is well-defined - i.e, there cannot be  $j \in S$  and  $j' \in L$  that both maximize  $\frac{1}{\alpha n} |\hat{C}_i \cap X_r|$  as  $r$  varies over  $[k]$ . To see this, WLOG, assume  $j \in S$ . Note that  $\frac{1}{\alpha n} |\hat{C}_i \cap X_j| \geq 1/k$ . Then, we immediately obtain:  $\frac{1}{\alpha n} |\cup_{j \in S} X_j \cap \hat{C}_i| \geq 1/k$ . Now, if we ensure that  $\eta' \leq \alpha^3/k^2$  and  $\varepsilon \leq O(\alpha^4/k)$ , then,  $\frac{1}{\alpha n} |\hat{C}_i \cap \cup_{j' \in L} X_{j'}|$  is at most the RHS of (4.7) which is  $\ll 1/k$ . This completes the proof of well-definedness. Next, adding up (4.7) for each  $i \in S$  yields that

$$\frac{1}{|\hat{C}_i|} |(\cup_{i \in G_S} X_i) \cap \cup_{j \in L} X_j| \leq O(\log(k/\eta')/\alpha) (\eta' + O(\varepsilon/\alpha)),$$

where we used that  $|G_S| \leq \ell$ . Combined with  $|\cup_{i \leq \ell} \hat{C}_i| \geq (1 - \eta'/k^{40})n$ , we obtain that

$$|\cup_{i \in G_S} X_i| \geq 1 - \eta'/k^{40} - O(\log(k/\eta')/\alpha) (\eta' + O(\varepsilon/\alpha)) = \eta + O(\log(k/\eta\alpha)\varepsilon/\alpha^2)$$

for  $\eta' \leq O(\eta^2\alpha^3/k)$ .

We now go ahead and establish (4.7). Let  $\tilde{\zeta}$  be a pseudo-distribution satisfying  $\mathcal{A}$  of degree  $(k/\eta)^{\text{poly}(k)}$  satisfying  $\tilde{\mathbb{E}}_{\tilde{\zeta}} z_i = \alpha$  for every  $i$ . Such a pseudo-distribution exists. To see why, let  $\tilde{\zeta}$  be the actual distribution that always sets  $X' = X$ , chooses an  $i$  with probability  $w_i$  and outputs a uniformly subset  $\hat{C}$  of size  $\alpha n$  of  $X_i$  conditioned on  $\hat{C}$  satisfying  $\mathcal{A}$ . Then, notice that since  $X$  satisfies Condition 2.47, by Fact 2.45, the uniform distribution on each  $X_i$  has  $t$ -certifiably  $C$ -hypercontractive degree 2 polynomials and is  $t$ -certifiably  $C\delta$ -anti-concentrated. By an concentration argument using high-order Chebyshev inequality, similar to the proof of Lemma 2.51 (applied to uniform distribution on  $X_i$  of size  $n \geq (dk)^{O(t)}$ ,  $\hat{C}$  chosen above satisfies the constraints  $\mathcal{A}$  with probability at least  $1 - o_d(1)$ . Observe that the probability that  $z_i$  is set to 1 under this distribution is then at most  $\alpha + o_d(1)$ . Thus, such a distribution satisfies all the constraints in  $\mathcal{A}$ .

Next, let  $M = \tilde{\mathbb{E}}_{\tilde{\zeta}}[zz^\top]$ . Then, we claim that:

1.  $o_d(1) + \alpha \geq M(i, j) \geq 0$  for all  $i, j$ ,
2.  $M(i, i) \in \alpha \pm o_d(1)$  for all  $i$ ,
3.  $\mathbb{E}_{j \sim [n]} M(i, j) \geq \alpha^2 - o_d(1)$  for every  $i$ .

The proofs of these basic observations are similar to those presented in Chapter 4.3 of [FKP19] (see also the proof of Theorem 5.1 in [BK20b]): Observe that  $\mathcal{A} \mid_{\frac{1}{4}} \{z_i z_j = z_i^2 z_j^2 \geq 0\}$  for every  $i, j$ . Thus, by Fact 2.18,  $\tilde{\mathbb{E}}[z_i z_j] \geq 0$  for every  $i, j$ . Next, observe that  $\mathcal{A} \mid_{\frac{1}{2}} \{(1 - z_i) = (1 - z_i)^2 \geq 0\}$  for every  $i$  and thus,  $\mathcal{A} \mid_{\frac{1}{2}} \{z_i(1 - z_j) \geq 0\}$ . Thus, by Fact 2.18 again, we must have  $\tilde{\mathbb{E}}[z_i z_j] \leq \tilde{\mathbb{E}}[z_i] \leq \alpha + o_d(1)$ . Finally,  $\mathcal{A} \mid_{\frac{1}{2}} \{\sum_j z_i z_j = z_i \sum_j z_j = \alpha n z_i\}$ . Thus, by Fact 2.18 again, we must have  $\sum_j M(i, j) = \sum_j \tilde{\mathbb{E}}[z_i z_j] = \alpha n \sum_j \tilde{\mathbb{E}}[z_j] \in (\alpha^2 \pm o_d(1))n$ . Let  $B_i$  be the entries in the  $i$ -th row  $M_i$  that are larger than  $\alpha^2/2$ . Then, by (1) and (2), we immediately derive that  $B_i$  must have at least  $\alpha n/2$  elements. Call an entry of  $M$  large if it exceeds  $\alpha^2 \eta'$ . For each  $i$ , let  $B_i$  be the set of large entries in row  $i$  of  $M$ . Then, using (3) and (1) above gives that  $|B_i| \geq \alpha(1 - \alpha \eta')n$  for each  $1 \leq i \leq n$ . Next, call a row  $i$  "good" if  $\frac{1}{\alpha n} \min\{|\cup_{r \in L} X_r \cap B_i|, |\cup_{r' \in S} X_{r'} \cap B_i|\} \leq 100\eta'/\alpha^3 + O(\varepsilon/\alpha^4)$ . Let us estimate the fraction of rows of  $M$  that are good.

Towards that goal, let us apply Lemma 4.5 with  $\eta$  set to  $\eta'$  and use Fact 2.18 (SoS Completeness), to obtain  $\sum_{r \in S, r' \in L} \mathbb{E}_{i \in X_r} \mathbb{E}_{j \in X_{r'}} M(i, j) \leq \eta' + O(\varepsilon/\alpha)$ . Using Markov's inequality, with probability  $1 - \alpha^3/100$  over the uniformly random choice of  $i, j$ ,  $\mathbb{E}_{j \in X_{r'}} M(i, j) \leq 100 \frac{1}{\alpha^3} \eta' + O(\varepsilon/\alpha^4)$ . Thus,  $1 - \alpha^3/100$  fraction of the rows of  $M$  are good.

Next, let  $R$  be the set of  $\frac{100}{\alpha} \log\left(\frac{k^{50}}{\eta'}\right)$  rows sampled in the run of the algorithm and set  $\hat{C}_i = B_i$  for every  $i \in R$ . The probability that all of them are good is then at least  $(1 - \alpha^3/100)^{\frac{100}{\alpha} \log\left(\frac{k^{50}}{\eta'}\right)} \geq 1 - \alpha$ . Let us estimate the probability that  $|\cup_{i \in R} \hat{C}_i| \geq (1 - \eta'/k^{40})n$ . For a uniformly random  $i$ , the chance that a given point  $t \in B_i$  is at least  $\alpha(1 - \alpha \eta')$ . Thus, the chance that  $t \notin \cup_{i \in R} B_i$  is at most

$(1 - \alpha/2)^{100/\alpha \log(k^{50}/(\alpha\eta'))} \leq \eta'/k^{50}$ . Thus, the expected number of  $t$  that are not covered by  $\cup_{i \in R} \hat{C}_i$  is at most  $n\eta'/k^{50}$ . Thus, by Markov's inequality, with probability at least  $1 - 1/k^{10}$ ,  $1 - \eta'/k^{40}$  fraction of  $t$  are covered in  $\cup_{i \in R} \hat{C}_i$ . By the above computations and a union bound, with probability at least  $1 - \eta'/k^{10}$  both the conditions below hold simultaneously: 1) each of the  $\frac{100}{\alpha} \log(k^{50}/\eta')$  rows  $R$  sampled are good and 2)  $|\cup_{i \in R} \hat{C}_i| \geq (1 - \eta'/k^{40})n$ . This completes the proof.  $\square$

### 4.3 Proof of Lemma 4.5

Our proof is based on the following simultaneous intersection bounds from [BK20b]. We will use the following lemma that forms the crux of the analysis of the clustering algorithm in [BK20b]:

**Lemma 4.7** (Simultaneous Intersection Bounds, Lemma 5.4 in [BK20b]). *Fix  $\delta > 0, k \in \mathbb{N}$ . Let  $X = X_1 \cup X_2 \cup \dots \cup X_k$  be a good sample of size  $n$  from a  $k$ -mixture  $\sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  of Gaussians. Let  $Y$  be any  $\varepsilon$ -corruption of  $X$ . Suppose there are  $r, r' \leq k$  such that one of the following three conditions hold for some  $\Delta \geq (k/\delta)^{O(k)}$ :*

1. *there exists a  $v$  such that  $v^\top \Sigma(r')v > \Delta v^\top \Sigma(r)v$  and  $B = \max_{i \leq k} \frac{v^\top \Sigma(i)v}{v^\top \Sigma(r')v}$ , or*
2. *there exists a  $v \in \mathbb{R}^d$  such that  $\langle \mu(r) - \mu(r'), v \rangle_2^2 \geq \Delta^2 v^\top (\Sigma(r) + \Sigma(r'))v$ , or,*
3.  $\|\Sigma(r')^{-1/2} \Sigma(r) \Sigma(r')^{-1/2} - I\|_F^2 \geq \Delta^2 \left( \|\Sigma(r')^{-1/2} \Sigma(r)^{1/2}\|_{\text{op}}^4 \right)$ .

Then, for the linear polynomial  $z(X_r) = \frac{1}{\alpha n} \sum_{i \in X_r} z_i$  in indeterminates  $z_i$ s satisfies:

$$\left\{ \cup_{i \leq 5} \mathcal{A}_i \mid \frac{z}{(k/\delta)^{O(k)} \log(2B)} \left\{ z(X_r)z(X_{r'}) \leq O(\sqrt{\delta}) + O(\varepsilon/\alpha) \right\} \right\}.$$

*Proof of Lemma 4.5.* Without loss of generality, assume that the pair of separated components are  $\mathcal{N}(\mu_1, \Sigma_1)$  and  $\mathcal{N}(\mu_2, \Sigma_2)$ . Let us start with the case when the pair is spectrally separated. Then, there is a  $v \in \mathbb{R}^d$  such that  $\Delta v^\top \Sigma_1 v \leq v^\top \Sigma_2 v$ .

Consider an ordering of the true clusters along the direction  $v$ , renaming cluster indices if needed, such that  $v^\top \Sigma_1 v \leq v^\top \Sigma_2 v \leq \dots \leq v^\top \Sigma_k v$ . Let  $j \leq k'$  be the largest integer such that  $\text{poly}(k/\eta) v^\top \Sigma_j v \leq v^\top \Sigma_{j+1} v$ . Further, observe that since  $j$  is defined to be the largest index which incurs separation  $\text{poly}(k/\eta)$ , all indices in  $[j, k]$  have spectral bound at most  $\text{poly}(k/\eta)$  and thus  $\frac{v^\top \Sigma_k v}{v^\top \Sigma_j v} \leq \text{poly}(k/\eta)^k$ .

Applying Lemma 4.7 with the above direction  $v$  to every  $r < j$  and  $r' \geq j$  and observing that the parameter  $B$  in each case is at most  $\frac{v^\top \Sigma_k v}{v^\top \Sigma_j v} \leq \Delta^k$  yields:

$$\mathcal{A} \mid \frac{z}{O(k^2 s^2 \text{poly} \log(\Delta))} \left\{ z(X_r)z(X_{r'}) \leq O(\varepsilon/\alpha) + \sqrt{\delta} \right\}.$$

Adding up the above inequalities over all  $r \leq j - 1$  and  $r' \geq j + 1$  and taking  $S = [j - 1]$ ,  $T = [k] \setminus [j - 1]$  completes the proof in this case.

Next, let us take the case when  $\mathcal{N}(\mu_1, \Sigma_1)$  and  $\mathcal{N}(\mu_2, \Sigma_2)$  are mean-separated. WLOG, suppose  $\langle \mu_1, v \rangle \leq \langle \mu_2, v \rangle \leq \dots \leq \langle \mu_k, v \rangle$ . Then, we know that  $\langle \mu_k - \mu_1, v \rangle \geq \Delta v^\top \Sigma_i v$ . Thus, there must exist

an  $i$  such that  $\langle \mu_i - \mu_{i+1}, v \rangle \geq \Delta v^\top \Sigma_i v / k$ . Let  $S = [i]$  and  $L = [k] \setminus S$ . Applying Lemma 4.7 and arguing as in the previous case (and noting that  $\kappa = \text{poly}(k)$ ) completes the proof.

Finally, let us work with the case of relative Frobenius separation. Since  $\|\Sigma_1^{-1/2} \Sigma_k^{1/2}\| \leq \text{poly}(k)$ , the hypothesis implies that  $\|\Sigma_1 - \Sigma_2\|_F \geq \Delta / \text{poly}(k)$ . Let  $B = \Sigma_1 - \Sigma_2$  and let  $A = B / \|B\|_F$ . WLOG, suppose  $\langle \Sigma_1, A \rangle \leq \dots \langle \Sigma_k, A \rangle$ . Then, since  $\langle \Sigma_k, A \rangle - \langle \Sigma_1, A \rangle \geq \Delta / \text{poly}(k)$ , there must exist an  $i$  such that  $\langle \Sigma_{i+1}, A \rangle - \langle \Sigma_i, A \rangle \geq \Delta / \text{poly}(k)$ . Let us now set  $S = [i]$  and  $L = [k] \setminus S$ .

Then, for every  $i \in S$  and  $j \in L$ , we must have:  $\langle \Sigma_j, A \rangle - \langle \Sigma_i, A \rangle \geq \Delta / \text{poly}(k)$ . Thus,  $\|\Sigma_j - \Sigma_i\|_F \geq \Delta / \text{poly}(k)$ . And thus,  $\Delta / \text{poly}(k) \leq \|\Sigma_j - \Sigma_i\|_F \leq \left\| \Sigma_i^{-1/2} \Sigma_j^{1/2} \right\|_2^2 \left\| \Sigma_i^{-1/2} \Sigma_j \Sigma_i^{-1/2} - I \right\|_F$ . Rearranging and using the bound on  $\left\| \Sigma_i^{-1/2} \Sigma_j^{1/2} \right\|_2^2$  yields that  $\left\| \Sigma_i^{-1/2} \Sigma_j \Sigma_i^{-1/2} - I \right\|_F \geq \Delta / \text{poly}(k)$ .

A similar argument as in the two cases above now completes the proof.  $\square$

#### 4.4 Proof of Lemma 4.6

We use  $\mathbb{E}_z$  as a shorthand for  $\frac{1}{\alpha n} \sum_{i=1}^n z_i$ . We will write  $\frac{1}{w_r n} \sum_{j \in X_r} z_j = z(X_r)$ . Note that  $z(X_r) \in [0, 1]$ . And finally, we will write  $z'(X_r) = \frac{1}{w_r n} \sum_{j \in X_r} z_j \mathbf{1}(x_j = y_j)$  – the version of  $z(X_r)$  that only sums over non-outliers.

We will use the following technical facts in the proof:

**Fact 4.8** (Lower Bounding Sums, Fact 4.19 [BK20b]). *Let  $A, B, C, D$  be scalar-valued indeterminates. Then, for any  $\tau > 0$ ,*

$$\{0 \leq A, B \leq A + B \leq 1\} \cup \{0 \leq C, D\} \cup \{C + D \geq \tau\} \Big|_{\frac{A, B, C, D}{2}} \{AC + BD \geq \tau AB\}.$$

**Fact 4.9** (Cancellation within SoS, Lemma 9.2 in [BK20b]). *For indeterminate  $a$  and any  $t \in \mathbb{N}$ ,*

$$\{a^{2t} \leq 1\} \Big|_{2t}^a \{a \leq 1\}.$$

**Lemma 4.10** (Lower-Bound on Variance of Degree 2 Polynomials). *Let  $Q \in \mathbb{R}^{d \times d}$ . Then, for any  $i, j \leq k$ , and  $z'(X_r) = \frac{1}{w_r n} \sum_{i \in X_r} z_i \mathbf{1}(x_i = y_i)$ , we have:*

$$\mathcal{A} \Big|_{\frac{z}{4}} \left\{ z'(X_r)^2 z'(X_r')^2 \leq \frac{(32Ct)^{2t}}{(\mathbb{E}_{X_r} Q - \mathbb{E}_{X_r'} Q)^{2t}} \left( \frac{\alpha^4}{w_r^2 w_r'^2} \left( \mathbb{E}_z (Q - \mathbb{E}_z Q)^2 \right)^t + \frac{\alpha^2}{w_r^2} \left( \mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^2 \right)^t + \frac{\alpha^2}{w_r'^2} \left( \mathbb{E}_{X_r'} (Q - \mathbb{E}_{X_r'} Q)^2 \right)^t \right) \right\}.$$

*Proof.* Let  $z'_i = z_i \mathbf{1}(x_i = y_i)$  for every  $i$ . Using the substitution rule and non-negativity constraints of the  $z_i$ 's, we have

$$\begin{aligned} \mathcal{A} \Big|_{\frac{z}{4}} \left\{ \mathbb{E}_z (Q - \mathbb{E}_z Q)^{2t} &= \frac{1}{\alpha^2 n^2} \sum_{i, j \leq n} z'_i z'_j \left( Q(x_i - x_j) - \mathbb{E}_z Q \right)^{2t} \\ &\geq \frac{1}{\alpha^2 n^2} \sum_{i, j \in X_r \text{ or } i, j \in X_r'} z'_i z'_j \left( Q(x_i - x_j) - \mathbb{E}_z Q \right)^{2t} \right\} \end{aligned} \quad (4.8)$$



Using the SoS almost triangle inequality, we have

$$\begin{aligned}
\mathcal{A} \Big|_4^z & \left\{ \frac{1}{\alpha^2 n^2} \sum_{i,j \in X_r \text{ or } i,j \in X_{r'}} z'_i z'_j \left( Q(x_i - x_j) - \mathbb{E}_z Q \right)^{2t} \right. \\
& \geq \left( \frac{1}{2^{2t}} \right) \left( \frac{1}{\alpha^2 n^2} \sum_{i,j \in X_r} z'_i z'_j \left( \mathbb{E}_{X_r} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha^2 n^2} \sum_{i,j \in X_r} z'_i z'_j \left( Q(x_i - x_j) - \mathbb{E}_{X_r} Q \right)^{2t} \right) \\
& \quad + \left( \frac{1}{2^{2t}} \right) \left( \frac{1}{\alpha^2 n^2} \sum_{i,j \in X_{r'}} z'_i z'_j \left( \mathbb{E}_{X_{r'}} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha^2 n^2} \sum_{i,j \in X_{r'}} z'_i z'_j \left( Q(x_i - x_j) - \mathbb{E}_{X_{r'}} Q \right)^{2t} \right) \quad (4.9) \\
& = 2^{-2t} \left( (w_r/\alpha)^2 z(X_r)^2 \left( \mathbb{E}_{X_r} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha^2 n^2} \sum_{i,j \in X_r} \left( Q(x_i - x_j) - \mathbb{E}_{X_r} Q \right)^{2t} \right) \\
& \quad + 2^{-2t} \left( (w_{r'}/\alpha)^2 z(X_{r'})^2 \left( \mathbb{E}_{X_{r'}} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha^2 n^2} \sum_{i,j \in X_{r'}} \left( Q(x_i - x_j) - \mathbb{E}_{X_{r'}} Q \right)^{2t} \right) \Big\}
\end{aligned}$$

Using Fact 4.8, we can further simplify the above as follows:

$$\begin{aligned}
\mathcal{A} \Big|_4^z & \left\{ \mathbb{E}_z (Q - \mathbb{E}_z Q)^{2t} \geq 2^{-6t} \frac{w_r^2 w_{r'}^2}{\alpha^4} z'(X_r)^2 z'(X_{r'})^2 \left( \mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q \right)^{2t} \right. \\
& \quad - 2^{-6t} (w_r/\alpha)^2 \mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^{2t} - 2^{-6t} (w_{r'}/\alpha)^2 \mathbb{E}_{X_{r'}} (Q - \mathbb{E}_{X_{r'}} Q)^{2t} \\
& \geq 2^{-6t} \frac{w_r^2 w_{r'}^2}{\alpha^4} z'(X_r)^2 z'(X_{r'})^2 \left( \mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q \right)^{2t} - (w_r/\alpha)^2 (Ct)^{2t} \left( \mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^2 \right)^t \\
& \quad \left. - (w_{r'}/\alpha)^2 (Ct)^{2t} \left( \mathbb{E}_{X_{r'}} (Q - \mathbb{E}_{X_{r'}} Q)^2 \right)^t \right\} \quad (4.10)
\end{aligned}$$

where the last inequality follows from the Certifiable Hypercontractivity constraint ( $\mathcal{A}_4$ ). Rearranging completes the proof.  $\square$

We can use the lemma above to obtain a simultaneous intersection bound guarantee when there are relative Frobenius separated components in the mixture.

**Lemma 4.11.** *Suppose  $\|\Sigma^{-1/2}(\Sigma_r - \Sigma_{r'})\Sigma^{-1/2}\|_F^2 \geq 10^8 \frac{C^6 t^4}{\beta^{2t} \alpha^2}$ . Then, for  $z'(X_r) = \frac{1}{\alpha n} \sum_{i \in X_r} z_i \cdot \mathbf{1}(y_i = x_i)$  for every  $r$ ,*

$$\mathcal{A} \Big|_{2t}^w \{ z'(X_r) z'(X_{r'}) \leq \beta \} .$$

*Proof.* We work with the transformation  $x_i \rightarrow \Sigma^{-1/2} x_i$ . Let  $\Sigma'_z = \Sigma^{-1/2} \Sigma_z \Sigma^{-1/2}$ ,  $\Sigma'_r = \Sigma^{-1/2} \Sigma_r \Sigma^{-1/2}$  and  $\Sigma'_{r'} = \Sigma^{-1/2} \Sigma_{r'} \Sigma^{-1/2}$  be the transformed covariances. Note that transformation is only for the purpose of the argument - our constraint system does not depend on  $\Sigma$ .

Notice that  $\|\Sigma'_{r'}\|_2 \leq \frac{1}{w_r}$  and  $\|\Sigma'_{r'}\|_2 \leq \frac{1}{w_{r'}}$ .

We now apply Lemma 4.10 with  $Q = \Sigma'_r - \Sigma'_{r'}$ . Then, notice that  $\mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q = \|\Sigma'_r - \Sigma'_{r'}\|_F^2 = \|Q\|_F^2$ . Then, we obtain:

$$\begin{aligned} & \mathcal{A} \Big|_{\frac{z}{2t}} \left\{ z'(X_r)^2 z'(X_{r'})^2 \right. \\ & \leq \left( \frac{32Ct}{\mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q} \right)^{2t} \left( \frac{\alpha^4}{w_r^2 w_{r'}^2} \left( \mathbb{E}_z (Q - \mathbb{E}_z Q)^2 \right)^t + \frac{\alpha^2}{w_r^2} \left( \mathbb{E}_{X_{r'}} (Q - \mathbb{E}_{X_{r'}} Q)^2 \right)^t + \frac{\alpha^2}{w_{r'}^2} \left( \mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^2 \right)^t \right) \Big\}. \end{aligned} \quad (4.11)$$

Since  $X_r$  and  $X_{r'}$  have certifiably  $C$ -bounded variance polynomials for  $C = 4$  (as a consequence of Condition 2.47 and Fact 2.45 followed by an application of Lemma 2.25), we have:

$$\mathcal{A} \Big|_{\frac{Q}{2}} \left\{ \mathbb{E}_{X_{r'}} (Q - \mathbb{E}_{X_{r'}} Q)^2 \leq 6 \left\| \Sigma'_{r'}{}^{1/2} Q \Sigma'_{r'}{}^{1/2} \right\|_F^2 \leq \frac{6}{w_{r'}^2} \|Q\|_F^2 \right\},$$

and

$$\mathcal{A} \Big|_{\frac{Q}{2}} \left\{ \mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^2 \leq 6 \left\| \Sigma_r{}^{1/2} Q \Sigma_r{}^{1/2} \right\|_F^2 \leq \frac{6}{w_r^2} \|Q\|_F^2 \right\}.$$

Finally, using the bounded-variance constraints in  $\mathcal{A}$ , we have:

$$\mathcal{A} \Big|_{\frac{Q,z}{4}} \mathbb{E} (Q - \mathbb{E}_z Q)^2 \leq \frac{6}{\alpha^2} \|Q\|_F^2.$$

Plugging these estimates back in (4.11) yields:

$$\mathcal{A} \Big|_{\frac{z}{4}} \left\{ z'(X_r)^2 z'(X_{r'})^2 \leq \frac{(1000Ct)^{2t}}{\alpha^{2t} \|Q\|_F^{2t}} \right\}. \quad (4.12)$$

Plugging in the lower bound on  $\|Q\|_F^{2t}$  and applying Fact 4.9 completes the proof.  $\square$

We can use the above lemma to complete the proof of Lemma 4.6:

*Proof of Lemma 4.6.* WLOG, assume that  $\Sigma = I$ . Let  $Q = \Sigma_r - \Sigma_{r'}$  and let  $\bar{Q} = Q/\|Q\|_F$ . Consider the numbers  $v_i = \text{tr}(\Sigma_r \cdot Q)$ . Then, we know that  $\max_{i,j} |v_i - v_j| \geq \|Q\|_F$ . Thus, there must exist a partition of  $[k] = S \cup T$  such that  $|v_i - v_j| \geq \|Q\|_F/k$  whenever  $i \in S$  and  $j \in T$ .

Thus, for every  $i \in S$  and  $j \in T$ ,  $\|\Sigma_i - \Sigma_j\|_F^2 \geq \|Q\|_F^2/k^2 = 10^8 \frac{C^6 t^4}{(\beta^{2t} \alpha^2)}$ . We can now apply Lemma above to every  $i \in S, j \in T$ , observe that  $\mathcal{A} \Big|_{\frac{z}{4}} \{z(X_r)z(X_{r'}) \leq z'(X_r)z'(X_{r'}) + 2\varepsilon/\alpha\}$ , and add up the resulting inequalities to finish the proof.  $\square$

## 4.5 Special Case: Algorithm for Uniform and Bounded Mixing Weights

In this subsection, we obtain a polynomial time algorithm when the input mixture has weights that are bounded from below. This includes the case of uniform weights and when the minimum mixing weight is at least some function of  $k$ . At a high level, our algorithm partitions the sample into clusters as long as there is a pair of components separated in TV distance and given samples that are not clusterable, runs the tensor decomposition algorithm to list decode. We then use standard robust tournament results to pick a hypothesis from the list.

**Theorem 4.12** (Robustly Learning Mixtures of Gaussians with Bounded Weights). *Given  $0 < \varepsilon < O_k(1)$ , let  $Y = \{y_1, y_2, \dots, y_n\}$  be a multiset of  $n \geq n_0 = \text{poly}_k(d, 1/\varepsilon)$   $\varepsilon$ -corrupted samples from a  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_{i \leq k} w_i \mathcal{N}(\mu_i, \Sigma_i)$ , such that  $w_i \geq \alpha$ . Then, there exists an algorithm with running time  $\text{poly}_k(n^{1/\alpha}) \cdot \exp(\text{poly}_k(1/\alpha, 1/\varepsilon))$  such that with probability at least  $9/10$  it outputs a hypothesis  $k$ -mixture of Gaussians  $\widehat{\mathcal{M}} = \sum_{i \leq k} \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  such that  $d_{TV}(\mathcal{M}, \widehat{\mathcal{M}}) = O_k(\varepsilon)$ .*

Briefly, our algorithm simply does the following:

1. **Clustering via SoS:** Guess a partition of the mixture such that each component in the partition is not clusterable. Let the resulting partition have  $t \leq k$  components. In parallel, try all possible ways to run Algorithm 4.4 repeatedly to obtain a partition of the samples,  $\{\tilde{Y}_j\}_{j \in [t]}$  into exactly  $t$  components. For each such partition repeat the following.
2. **Robust Isotropic Transformation:** Run the algorithm corresponding to Lemma 6.7 on each set  $\tilde{Y}_j$  to make the sample approximately isotropic. Grid search for weights over  $[\alpha, 1/k]^k$  with precision  $\alpha$ .
3. **List-Decoding via Tensor Decomposition:** Run Algorithm 3.2 on each  $\tilde{Y}_j$ . Concatenate the lists to obtain  $\mathcal{L}$ .
4. **Robust Tournament:** Run the tournament from Fact 2.52 over all the hypotheses in  $\mathcal{L}$ , and output the winning hypothesis.

*Proof Sketch.* Setting  $\Delta = (k^{k^{O(k)}})$ , it follows from Theorem 4.2 that we obtain a partition of  $Y$  into  $\{\tilde{Y}_j\}_{j \in [t]}$ , for some  $t \in [k]$  such that  $\tilde{Y}_j$  has at most  $O(k\varepsilon/\alpha)$  outliers,  $(1 - O(k\varepsilon/\alpha))$ -fraction of samples from at least one component of the input mixture and the resulting samples are not  $\Delta$ -separated (see Definition 4.1). It then follows from Lemma 6.7 that the mean  $\mu_j$  and covariance  $\Sigma_j$  of  $\tilde{Y}_j$  satisfy : a)  $\|\mu_j\|_2 \leq O(\sqrt{\varepsilon}k^{1.5}/\alpha^{1.5})$ , b)  $(1 - \sqrt{\varepsilon}k^{1.5}/\alpha^{1.5})I \leq \Sigma_j \leq (1 + \sqrt{\varepsilon}k^{1.5}/\alpha^{1.5})I$ , and c)  $\|\Sigma_j - I\|_F \leq O(\sqrt{\varepsilon}k^{1.5}/\alpha^{1.5})$ .

Each component,  $\tilde{Y}_j$ , of the partition can have at most  $k$  components. Assuming these correspond to  $\{w_i^{(j)}, \mu_i^{(j)}, \Sigma_i^{(j)}\}_{i \in [k]}$ , observe,  $\sum_{i \in [k]} w_i^{(j)} \Sigma_i^{(j)} + w_i^{(j)} \mu_i^{(j)} (\mu_i^{(j)})^\top \leq (1 + \sqrt{\varepsilon}k^{1.5}/\alpha^{1.5})I$ . Thus, we have that  $\|\mu_i^{(j)}\|_2^2 \leq (1 + \sqrt{\varepsilon}k^{1.5})/\alpha^{2.5}$  and combined with not being  $\Delta$ -separated, it follows that for all

$i' \in [k]$ ,

$$\begin{aligned}
\left\| \Sigma_{i'}^{(j)} - I \right\|_F &= \left\| \Sigma_{i'}^{(j)} - \Sigma_j + (\Sigma_j - I) \right\|_F \leq \left\| \Sigma_{i'}^{(j)} - \sum_{i \in [k]} w_i^{(j)} \Sigma_i^{(j)} + \sum_{i \in [k]} w_i^{(j)} \mu_i^{(j)} \left( \mu_i^{(j)} \right)^\top \right\|_F + \left\| \Sigma_j - I \right\|_F \\
&\leq \left\| \sum_{i \in [k]} w_i^{(j)} \left( \Sigma_{i'}^{(j)} - \Sigma_i^{(j)} \right) \right\|_F + \mathcal{O}(k^{1.5}/\alpha^{2.5}) \\
&\leq \mathcal{O}(\Delta/\alpha) .
\end{aligned}$$

There are at most  $\mathcal{O}(k^k)$  ways in which we can partition the set of input points such that each resulting component is not partially clusterable. We run the algorithm in parallel for each one. Then, for the correct iteration, we apply Theorem 3.1 to get a list  $\mathcal{L}$  of size  $\exp(\text{poly}_k(1/\alpha, 1/\varepsilon))$  such that it contains a hypothesis  $\{\hat{w}_i^{(j)}, \hat{\mu}_i^{(j)}, \hat{\Sigma}_i^{(j)}\}_{i \in [k]}$  such that  $|\hat{w}_i^{(j)} - w_i^{(j)}| \leq \alpha$ ,  $\|\hat{\mu}_i^{(j)} - \mu_i^{(j)}\|_2 \leq \mathcal{O}_k(\varepsilon)$  and  $\|\hat{\Sigma}_i^{(j)} - \Sigma_i^{(j)}\|_F \leq \mathcal{O}_k(\varepsilon)$ . Since  $(1 - 1/\Delta)I \leq \Sigma_i^{(j)}$ , it then follows from Lemma 6.5 that the hypothesis is  $\mathcal{O}_k(\varepsilon)$ -close to the input in total variation distance.

Algorithm 4.4 is called at most  $\mathcal{O}(k^k)$  times, and along with the robust isotropic transformation, this requires  $\text{poly}_k(n^{1/\alpha}, 1/\varepsilon)$ . The grid search contributes a multiplicative factor of  $(1/\alpha)^k$ . The tensor decomposition algorithm and robust hypothesis section  $\text{poly}_k(n^{1/\alpha}) \cdot \exp(\text{poly}_k(1/\alpha, 1/\varepsilon))$ .  $\square$

## 5 Spectral Separation of Thin Components

In this section, we show how to efficiently separate a thin component, if such a component exists, given sufficiently accurate approximations to the component means and covariances. This is an important step in our overall algorithm and is required to obtain total variation distance guarantees.

Specifically, the main algorithmic result of this section is described in the following lemma:

**Lemma 5.1.** *There is a polynomial-time algorithm with the following properties: Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians on  $\mathbb{R}^d$ , and let  $X$  be a set of points in  $\mathbb{R}^d$  satisfying Condition 2.47 with respect to  $\mathcal{M}$  for some parameters  $(\gamma, t)$ . The algorithm takes input parameters  $\eta, \delta$ , satisfying  $0 < \delta < \eta < 1/(100k)$ , and  $Y$ , an  $\varepsilon$ -corrupted version of  $X$ , as well as candidate parameters  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \leq k}$ . Then as long as*

1.  $\text{Cov}(\mathcal{M}) \geq I/2$ ,
2.  $\|\mu_i - \hat{\mu}_i\|_2 < \delta$  and  $\|\Sigma_i - \hat{\Sigma}_i\|_F < \delta$ , for all  $i \in [k]$ , and
3. there exists an  $s \in [k]$  such that  $\Sigma_s$  has an eigenvalue  $< \eta$ ,

*the algorithm outputs a partition of  $Y$  into  $Y_1 \cup Y_2$  such that there is a non-trivial partition of  $[k]$  into  $Q_1 \cup Q_2$ , so that letting  $\mathcal{M}_j$ ,  $j \in \{1, 2\}$ , be proportional to  $\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then  $Y_j$  is an  $((\mathcal{O}(k^2\gamma) + \tilde{\mathcal{O}}(\eta^{1/2k}))/W_j)$ -corruption of a set satisfying Condition 2.47 with respect to  $\mathcal{M}_j$  with parameters  $(\mathcal{O}(k\gamma/W_j), t)$ .*

The key component in the proof of Lemma 5.1 is the following lemma:

**Lemma 5.2.** Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians in  $\mathbb{R}^d$  with  $\text{Cov}(\mathcal{M}) \geq I/2$ . Suppose that, for some  $0 < \delta < 1/(100k)$ , we are given  $\hat{\mu}_i$  and  $\hat{\Sigma}_i$  satisfying  $\|\mu_i - \hat{\mu}_i\|_2 < \delta$  and  $\|\Sigma_i - \hat{\Sigma}_i\|_F < \delta$ , for all  $i \in [k]$ . Suppose furthermore that for some  $\eta > \delta$ , there is a  $\Sigma_s$ ,  $s \in [k]$ , with an eigenvalue less than  $\eta$ . There exists a computationally efficient algorithm that takes inputs  $\eta, \delta, \hat{\mu}_i, \hat{\Sigma}_i$ , and computes a function  $F : \mathbb{R}^d \rightarrow \{0, 1\}$  such that:

1. For each  $i \in [k]$ ,  $F(G_i)$  returns the same value in  $\{0, 1\}$  with probability at least  $1 - \tilde{O}_k(\eta^{1/(2k)})$ . We define the most likely value of  $F(G_i)$  to be this value.
2. There exist  $i, j \in [k]$  such that the most likely values of  $F(G_i)$  and  $F(G_j)$  are different.

Furthermore,  $F(x)$  can be chosen to be of the form  $f(v \cdot x)$ , for some  $v \in \mathbb{R}^d$ , and  $f : \mathbb{R} \rightarrow \{0, 1\}$  is an  $O(k)$ -piecewise constant function.

Given Lemma 5.2, it is easy to finish the proof of Lemma 5.1.

*Proof of Lemma 5.1.* We simply take the candidate parameters, obtain  $F$  from Lemma 5.2, and partition  $Y = Y_1 \cup Y_2$ , so that  $F$  is constant on both  $Y_1$  and  $Y_2$ . We let  $Q_j$  be the set of  $i$  so that  $F(G_i)$  returns the value  $j - 1$  with large probability. Letting the partition of  $X$  for Condition 2.47 be  $X = X_1 \cup \dots \cup X_k$ , we let  $X^j = \bigcup_{i \in Q_j} X_i$ . Lemma 2.50 shows that the  $X^j$  satisfy the appropriate conditions for  $\mathcal{M}_j$ . It remains to prove that  $Y_j$  equals  $X^j$  with a sufficiently small rate of corruptions. The fraction of points misclassified by  $F$  equals  $\varepsilon$  (the fraction of outliers in the sample  $Y$ ) plus the misclassification error of  $F$ . We note that given the form of  $F$  and the fact that the uncorrupted samples in  $Y$  satisfy Condition 2.47, the fraction of misclassified samples from each component  $i$  is at most the probability that a random sample from  $G_i$  gets misclassified (at most  $\tilde{O}_k(\eta^{1/(2k)})$  by Lemma 5.2) plus  $O(k\gamma)$ . Summing this over components, gives Lemma 5.1.  $\square$

Let us now describe the algorithm to prove Lemma 5.2 (and evaluate  $F$ ), which is given in pseudocode below (Algorithm 5.3).

**Algorithm 5.3** (Algorithm for Spectrally Separating Thin Components).

**Input:** Estimated parameters  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \leq k}$ , parameters  $\eta, \delta$ .

**Output:** A function  $F : \mathbb{R}^d \rightarrow \{0, 1\}$ .

**Operation:**

1. Find a unit-norm direction  $v$  such that there exists  $s \in [k]$ ,  $v^T \hat{\Sigma}_s v < 2\eta$ .
2. Compute  $(v^T \hat{\Sigma}_i v)$  for all  $i \in [k]$ .
  - (a) If there exists  $j \in [k]$  such that  $(v^T \hat{\Sigma}_j v) > \sqrt{\eta}$ , find a  $t$  such that  $\sqrt{\eta} > t > 2\eta$  and there is no  $j \in [k]$  with  $t < v^T \hat{\Sigma}_j v < t \Omega(\eta^{-1/(2k)})$ . Set  $F(x) = 1$  if there is an  $i$  such that  $|v \cdot (x - \hat{\mu}_i)| < \sqrt{t} \log(1/\eta)$  and 0 otherwise.

(b) Otherwise, compute  $v \cdot \hat{\mu}_i$  for all  $i \in [k]$ . Find a  $t$  between the minimum and the maximum of  $v \cdot \hat{\mu}_i$  such that there is no  $v \cdot \hat{\mu}_i$  within  $1/(20k)$  of  $t$ . Set  $F(x) = 1$  if  $v \cdot x > t$  and 0 otherwise.

*Proof of Lemma 5.2.* Let  $v$  be a unit vector and  $s \in [k]$  such that  $v^T \hat{\Sigma}_s v < 2\eta$ . By assumption, we have that  $\mathbf{Var}[v \cdot \mathcal{M}] \geq 1/2$ . Furthermore,

$$\mathbf{Var}[v \cdot \mathcal{M}] = \sum_i w_i (v^T \Sigma_i v) + \sum_i w_i (v \cdot (\mu_i - \mu))^2 \leq \sum_i w_i (v^T \Sigma_i v) + \sum w_i (v \cdot (\mu_i - \mu_s))^2,$$

where  $\mu$  is the mean of  $\mathcal{M}$ . This means that either there exists  $j \in [k]$  such that  $(v^T \Sigma_j v) > 1/4$ , or there exists  $j \in [k]$  such that  $|v \cdot (\mu_j - \mu_s)| > 1/4$ . Since we have approximations of these quantities to order  $\delta$ , we have that there is  $j \in [k]$  such that  $(v^T \hat{\Sigma}_j v) > 1/10$  or that there is  $j \in [k]$  with  $|v \cdot (\hat{\mu}_j - \hat{\mu}_s)| > 1/10$ .

We first consider the case that there is a  $j \in [k]$  such that  $(v^T \hat{\Sigma}_j v) > \sqrt{\eta}$ . Since there is a  $j \in [k]$  with  $(v^T \hat{\Sigma}_j v) > \sqrt{\eta}$  and another  $s \in [k]$  with  $(v^T \hat{\Sigma}_s v) < 2\eta$ , there must be some  $\sqrt{\eta} > t > 2\eta$  such that there is no  $j \in [k]$  with  $t < v^T \hat{\Sigma}_j v < t \Omega(\eta^{-1/(2k)})$ . Otherwise, there must be at least one  $\hat{\Sigma}_i$  in each  $2\eta \leq \Omega(\eta^{-1/(2k)})^i \leq \sqrt{\eta}$ , where we need more than  $k$  components.

For a given  $x$ , we define  $F(x)$  to be 1 if there exists  $i$  such that  $|v \cdot (x - \hat{\mu}_i)| < \sqrt{t} \log(1/\eta)$ , and  $F(x) = 0$  otherwise.

To show that this works, we note that for all  $i \in [k]$ , if  $v^T \hat{\Sigma}_i v \leq t$ , then  $\mathbf{Var}[v \cdot G_i] \leq t + \delta$ , and since  $|v \cdot (\mu_i - \hat{\mu}_i)| < \delta$ , by the Gaussian tail bound, we have that

$$\mathbb{P}_{x \sim G_i} \left( |x - \mu_i| \geq (\sqrt{t} \log(1/\eta) - \delta) \right) \leq \exp \left( -\frac{(\sqrt{t} \log(1/\eta) - \delta)^2}{2(t + \delta)} \right) = O(\eta).$$

Thus, all but an  $O(\eta)$ -fraction of the samples of  $G_i$  have  $F(x) = 1$ .

On the other hand, for components  $i$  with  $v^T \hat{\Sigma}_i v \gg t\eta^{-1/(2k)}$ , we have that  $\mathbf{Var}[v \cdot G_i] \gg t\eta^{-1/(2k)}$ . Then, the density of  $G_i$  is at most  $1/\sqrt{2\pi t\eta^{-1/(2k)}}$ . So, the probability that a sample from  $v \cdot G_i$  lies in any interval of length  $2\sqrt{t} \log(1/\eta)$  is at most

$$\frac{1}{\sqrt{2\pi t\eta^{-1/(2k)}}} 2\sqrt{t} \log(1/\eta) = \tilde{O}(\eta^{1/(4k)}).$$

Since there are  $k$  such intervals, the probability that  $F(x)$  is 1 when  $x$  is drawn from  $G_i$  is at most  $\tilde{O}_k(\eta^{1/(4k)})$ . This completes our proof of point (1), and point (2) follows from the fact that we know of component  $G_j$  in one class and  $G_s$  in the other class.

We next consider the case where  $(v^T \hat{\Sigma}_j v) \leq \sqrt{\eta}$  for all  $j \in [k]$ , and where  $|v \cdot (\hat{\mu}_j - \hat{\mu}_s)| > 1/10$  for some  $j \in [k]$ . Then we can find some  $t$  between  $v \cdot \hat{\mu}_j$  and  $v \cdot \hat{\mu}_s$  such that no  $v \cdot \hat{\mu}_i$  is within  $1/(20k)$  of  $t$ . In this case, we define  $F(x)$  be 1 if  $v \cdot x > t$  and 0 otherwise. To show part (1), first consider  $i \in [k]$  such that  $v \cdot \hat{\mu}_i < t - 1/(20k)$ . Then we have that  $v \cdot \mu_i < t - 1/(30k)$ . Furthermore,  $\mathbf{Var}[v \cdot G_j] \leq \delta + \sqrt{\eta}$ . Therefore, the probability that  $v \cdot G_i > t$  is at most  $\exp(-\Omega_k((\delta + \sqrt{\eta})^{-2}))$ , which is sufficient.

A similar argument holds in the other direction for  $i \in [k]$  such that  $v \cdot \hat{\mu}_i > t + 1/(20k)$ , and statement (2) holds because we know that there are both kinds of components. This completes the proof.  $\square$

## 6 Robust Proper Learning: Proof of Theorem 1.4

In this section, we show how to combine the partial clustering, tensor decomposition, and recursive clustering algorithms to establish our main result. The main theorem we prove is as follows:

**Theorem 6.1** (Robustly Learning  $k$ -Mixtures of Arbitrary Gaussians). *Given  $0 < \varepsilon < 1/k^{k^{O(k^2)}}$  and a multiset  $Y = \{y_1, y_2, \dots, y_n\}$  of  $n$  i.i.d. samples from a distribution  $F$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , for an unknown  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_{i \leq k} w_i \mathcal{N}(\mu_i, \Sigma_i)$ , where  $n \geq n_0 = d^{O(k)}/\text{poly}(\varepsilon)$ , Algorithm 6.3 runs in time  $n^{O(1)} \exp(O(k)/\varepsilon^2)$  and with probability at least 0.99 outputs a hypothesis  $k$ -mixture of Gaussians  $\widehat{\mathcal{M}} = \sum_{i \leq k} \hat{w}_i \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i)$  such that  $d_{\text{TV}}(\mathcal{M}, \widehat{\mathcal{M}}) = O(\varepsilon^{c_k})$ , with  $c_k = 1/(100^k C^{(k+1)!} k! \text{sf}(k+1))$ , where  $C > 0$  is a universal constant and  $\text{sf}(k) = \prod_{i \in [k]} (k-i)!$  is the super-factorial function.*

As an immediate corollary, we obtain the following:

**Corollary 6.2** (Robustly Learning  $k$ -Mixtures of Gaussians in Polynomial Time). *Given  $0 < \varepsilon < 1/\exp(k^{k^{O(k^2)}})$ , and a multiset  $Y = \{y_1, y_2, \dots, y_n\}$  of  $n$  i.i.d. samples from a distribution  $F$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , for an unknown  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_{i \leq k} w_i \mathcal{N}(\mu_i, \Sigma_i)$ , where  $n \geq n_0 = d^{O(k)} \log^{O(1)}(1/\varepsilon)$ , there exists an algorithm that runs in time  $\text{poly}_k(n, 1/\varepsilon)$  and with probability at least 0.99 outputs a  $k$ -mixture of Gaussians  $\widehat{\mathcal{M}} = \sum_{i \leq k} \hat{w}_i \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i)$  such that  $d_{\text{TV}}(\mathcal{M}, \widehat{\mathcal{M}}) = O\left((1/\log(1/\varepsilon))^{1/(k^{O(k^2)})}\right)$ .*

The corollary follows by running Algorithm 6.3 with  $\varepsilon \leftarrow \sqrt{1/\log(1/\varepsilon)}$  and applying Theorem 6.1.

The algorithm establishing Theorem 6.1 is given in pseudocode below. Algorithm 6.4 takes as input a corrupted sample from a  $k$ -mixture of Gaussians and outputs a set of  $k$  mixing weights, means, and covariances, such that the resulting mixture is close to the input mixture in total variation distance with non-negligible probability. Algorithm 6.3 simply runs Algorithm 6.4 many times to create a small list of candidate hypotheses (consisting of mixing weights, means, and covariances), and finally runs a robust tournament to outputs a winner. This boosts the probability of success to at least 0.99.

**Algorithm 6.3** (Algorithm for Robustly Learning Arbitrary GMMs).

**Input:** An outlier parameter  $\varepsilon > 0$  and a component-number parameter  $k \in \mathbb{N}$ . An  $\varepsilon$ -corrupted sample  $Y = \{y_1, y_2, \dots, y_n\}$  from a  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_{i \in [k]} w_i \mathcal{N}(\mu_i, \Sigma_i)$ .

**Parameters:** Let  $c_k = 1/(100^k C^{(k+1)!} \text{sf}(k+1)k!)$  be a scalar function of  $k$ , where  $\text{sf}(k) = \prod_{i \in [k]} (k-i)!$  and  $C$  is a sufficiently large constant.



**Output:** A set of parameters  $\{(\hat{w}_i, \hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]}$ , such that with probability at least 0.99 the mixture  $\widehat{\mathcal{M}} = \sum_{i \in [k]} \hat{w}_i \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i)$  is  $\mathcal{O}(\varepsilon^{c_k})$ -close in total variation distance to  $\mathcal{M}$ .

**Operation:**

1. Let  $\mathcal{L} = \{\phi\}$  be an empty list. Repeat the following  $\exp(O(k)/\varepsilon^2)$  times :
  - (a) Run Algorithm 6.4 with input  $Y$ , fraction of outliers  $\varepsilon$ , and number of components  $k$ . Let the resulting output be a set of  $k$  mixing weights, means and covariances, denoted by  $\{(\hat{w}_i, \hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]}$ . Add  $\{(\hat{w}_i, \hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]}$  to  $\mathcal{L}$ .
2. Run the robust tournament from Fact 2.52 over all the hypotheses in  $\mathcal{L}$ . Output the winning hypothesis, denoted by  $\{(\hat{w}_i, \hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]}$ .

**Algorithm 6.4** (Cluster or List-Decode).

**Input:** An outlier parameter  $0 < \varepsilon < 1$  and a component-number parameter  $k \in \mathbb{N}$ . An  $\varepsilon$ -corrupted version  $Y = \{y_1, y_2, \dots, y_n\}$  of  $X$ , where  $X$  is a set of  $n$  samples from a  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_{i \in [k]} w_i \mathcal{N}(\mu_i, \Sigma_i)$  such that  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\varepsilon d^{-8k} k^{-C'k}, 8k + 48)$ , where  $C' > 0$  is a sufficiently large constant.

**Parameters:** Let  $c_k = 1/(100^k C^{(k+1)!} \text{sf}(k+1)k!)$  be a scalar function of  $k$ , where  $\text{sf}(k) = \prod_{i \in [k]} (k-1)!$  and  $C$  is a sufficiently large constant.

**Output:** A set of parameters  $\{(\hat{w}_i, \hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]}$  such that with probability at least  $\exp(-O(k)/\varepsilon^2)$ ,  $d_{TV}\left(\sum_{i \in [k]} \hat{w}_i \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i), \mathcal{M}\right) \leq \mathcal{O}(\varepsilon^{c_k})$ .

**Operation:**

1. **Treat Light Component as Noise:** If  $k = 0$ , ABORT. With probability 1/2, run Algorithm 6.4 on samples  $Y$ , with fraction of outliers  $\varepsilon + \varepsilon^{1/(10C^{k+1}(k+1)!)}$  and number of components  $k - 1$ . Return the resulting set of estimated parameters,  $\{(\hat{w}_i, \hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k-1]}$ , appended with  $(0, 0, I)$ . Else, do the following:
  - // We guess whether the event that the minimum mixing weight  $\alpha$  is at least  $\varepsilon^{1/(10C^{k+1}(k+1)!)}$
  - // holds. If it does not, we proceed with the algorithm. Else, we treat the smallest weight
  - // component as noise and recurse with  $k - 1$  components.
2. **Robust Isotropic Transformation:** With probability 0.5, run the algorithm corresponding to Lemma 6.7 on the samples  $Y$ , and let  $\hat{\mu}, \hat{\Sigma}$  be the robust estimates of the mean and covariance. If  $k = 1$ , return  $(\hat{w} = 1, \hat{\mu}, \hat{\Sigma})$ . Else, compute  $\hat{U} \hat{\Lambda} \hat{U}^\top$ , the eigendecomposition of  $\Sigma$ , and for all  $i \in [n]$ , apply the affine transformation  $y_i \rightarrow \hat{U}^\top \hat{\Sigma}^{+1/2} (y_i - \hat{\mu})$ .
  - // The resulting estimates  $\hat{\mu}, \hat{\Sigma}$  satisfy Lemma 6.7, and the uncorrupted samples are
  - // effectively drawn from a nearly isotropic  $k$ -mixture.

3. With probability  $1/2$ , run either (a) or (b) in the following:

(a) **Partial Clustering via SoS:** Run Algorithm 4.4 with outlier parameter  $\varepsilon$  and accuracy parameter  $\varepsilon^{1/(5C^{k+1}(k+1)!)}$ . Let  $Y_1, Y_2$  be the partition returned. Guess the number of components in  $Y_1$  to be some  $k_1 \in [k-1]$  uniformly at random. Run Algorithm 6.4 with input  $Y_1$ , fraction of outliers  $\varepsilon^{1/(10c^{k+1}(k+1)!)}$ , and number of components  $k_1$ , and let  $\left\{(\hat{w}_i^{(1)}, \hat{\mu}_i^{(1)}, \hat{\Sigma}_i^{(1)})\right\}_{i \in [k_1]}$  be the resulting output. Similarly, run Algorithm

6.4 with input  $Y_2$ , fraction of outliers  $\varepsilon^{1/(10c^{k+1}(k+1)!)}$ , and number of components  $k - k_1$ , and let  $\left\{(\hat{w}_i^{(2)}, \hat{\mu}_i^{(2)}, \hat{\Sigma}_i^{(2)})\right\}_{i \in [k-k_1]}$  be the resulting output. Output the set

$$\left\{(\hat{w}_i^{(1)}|Y_1|/|Y|, \hat{\mu}_i^{(1)}, \hat{\Sigma}_i^{(1)})\right\}_{i \in [k_1]} \cup \left\{(\hat{w}_i^{(2)}|Y_2|/|Y|, \hat{\mu}_i^{(2)}, \hat{\Sigma}_i^{(2)})\right\}_{i \in [k-k_1]}.$$

// When the mixture is covariance separated, the preconditions of Theorem 4.3 are

// satisfied (see Lemma 6.8). The partition is non-trivial, and the fraction of outliers

// increases from  $\varepsilon \rightarrow \varepsilon^{1/(10c^{k+1}(k+1)!)}$ .

(b) **List-Decoding via Tensor Decomposition:** Run Algorithm 3.2 and let  $L$  be the resulting list of hypotheses such that each hypothesis is a set of parameters  $\{(\hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]}$ . Let  $\tau = \Theta\left(\varepsilon^{1/(40C^{k+1}(k+1)!)}\right)$  be an eigenvalue threshold. Select a hypothesis,  $\{(\hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]} \in L$  uniformly at random.

// Conditioned on not being covariance separated, we satisfy the preconditions of

// Theorem 3.1 (see Lemma 6.9). The output is a list that contains  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$

// such that for all  $i \in [k]$ ,  $\|\hat{\mu}_i - \mu_i\|_2 = \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$  and

//  $\left\|\hat{\Sigma}_i - \Sigma_i\right\|_F = \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$ .

i. **Large Eigenvalues:** If for all  $i \in [k]$ ,  $\hat{\Sigma}_i \geq \tau I$ , sample  $\hat{w}_i$  from  $[0, 1]$  uniformly at random such that  $\sum_i \hat{w}_i = (1 \pm k\varepsilon)$ . Return  $\left\{(\hat{w}_i, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \hat{\mu}_i + \hat{\mu}, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \hat{\Sigma}_i \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top)\right\}_{i \in [k]}$ .

// If all estimated covariances have all eigenvalues larger than  $\tau$ , the recursion // bottoms out and the hypothesis is returned.

ii. **Spectral Separation of Thin Components:** Else,  $\exists v, i$  s.t.  $v^\top \hat{\Sigma}_i v \leq \tau$ . Run the algorithm corresponding to Lemma 5.1 with input  $Y$ , parameter estimates  $\{(\hat{\mu}_i, \hat{\Sigma}_i)\}_{i \in [k]}$  and threshold  $\tau$ . Let  $Y_1$  and  $Y_2$  be the resulting partition.

// Use small eigenvalue directions to partition the points.

A. If  $\min(|Y_1|, |Y_2|) < \varepsilon^{1/(400kC^{k+1}(k+1)!)}n$ , run Algorithm 6.4 with input  $Y$ , fraction of outliers  $2\varepsilon^{1/(400kC^{k+1}(k+1)!)}$  and number of components being  $k-1$ , and let  $\left\{(\hat{w}_i^{(1)}, \hat{\mu}_i^{(1)}, \hat{\Sigma}_i^{(1)})\right\}_{i \in [k_1]}$  be the resulting output. Output the resulting hypothesis

$$\left\{(\hat{w}_i, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \hat{\mu}_i + \hat{\mu}, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \hat{\Sigma}_i \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top)\right\}_{i \in [k-1]} \cup (0, 0, I).$$

B. Else, select  $k_1 \in [k-1]$  uniformly at random. Run Algorithm 6.4 with input  $Y_1$ , fraction of outliers  $\varepsilon^{1/(100kC^{k+1}(k+1)!)}$  and number of

components being  $k_1$ . Similarly, run Algorithm 6.4 with input  $Y_2$ , fraction of outliers  $\varepsilon^{1/(100kC^{k+1}(k+1)!)}$  and number of components  $k - k_1$ , and let  $\{(\hat{w}_i^{(2)}, \hat{\mu}_i^{(2)}, \hat{\Sigma}_i^{(2)})\}_{i \in [k-k_1]}$  be the resulting output. Output the set  $\{(\hat{w}_i^{(1)}|Y_1|/|Y|, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \hat{\mu}_i^{(1)} + \hat{\mu}, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \hat{\Sigma}_i^{(1)}\hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top)\}_{i \in [k_1]} \cup \{(\hat{w}_i^{(2)}|Y_2|/|Y|, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \mu_i^{(2)} + \hat{\mu}, \hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top \hat{\Sigma}_i^{(2)}\hat{U}\hat{\Lambda}^{1/2}\hat{U}^\top)\}_{i \in [k-k_1]}$ .

## 6.1 Analysis of Algorithm 6.3

To prove Theorem 6.1, we will require the following intermediate results. We defer some proofs in this subsection to Appendix A.

We use the following lemma to relate the Frobenius distance of covariances to the total variation distance between two Gaussians, when the eigenvalues of the covariances are bounded below.

**Lemma 6.5** (Frobenius Distance to TV Distance). *Suppose  $\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)$  are Gaussians with  $\|\mu_1 - \mu_2\|_2 \leq \delta$  and  $\|\Sigma_1 - \Sigma_2\|_F \leq \delta$ . If the eigenvalues of  $\Sigma_1$  and  $\Sigma_2$  are at least  $\lambda > 0$ , then  $d_{\text{TV}}(\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)) = O(\delta/\lambda)$ .*

We start by showing that when Condition 2.47 holds, the uniform distribution on a  $(1 - \varepsilon)$ -fraction of the points is certifiably hypercontractive.

**Lemma 6.6** (Component Moments to Mixture Moments). *Let  $\mathcal{M} = \sum_{i \in [k]} w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture with mean  $\mu$  and covariance  $\Sigma$  such that  $w_i \geq \alpha$ , for some  $0 < \alpha < 1$ , and for all  $i, j \in [k]$ ,  $\|\Sigma^{+1/2}(\Sigma_i - \Sigma_j)\Sigma^{+1/2}\|_F \leq 1/\sqrt{\alpha}$ . Let  $X$  be a multiset of  $n$  samples satisfying Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$ , for  $0 < \gamma < (dk/\alpha)^{-ct}$ , for a sufficiently large constant  $c$ , and  $t \in \mathbb{N}$ . Let  $\mathcal{D}$  be the uniform distribution over  $X$ . Then,  $\mathcal{D}$  is  $2t$ -certifiably  $(c/\alpha)$ -hypercontractive and for  $d \times d$ -matrix-valued indeterminate  $Q$ ,  $\frac{Q}{2} \left\{ \mathbb{E}_{\mathcal{M}}(x^\top Q x - \mathbb{E}_{\mathcal{M}} x^\top Q x)^2 \leq O(1/\alpha) \|\Sigma^{1/2} Q \Sigma^{1/2}\|_F^2 \right\}$ .*

Next, we show how to robustly estimate the mean and covariance of an  $\varepsilon$ -corrupted set of samples satisfying Condition 2.47 when the mixture is not partially clusterable, and make the inliers nearly isotropic.

**Lemma 6.7** (Robust Isotropic Transformation). *Given  $0 < \varepsilon < 1$ , and  $k \in \mathbb{N}$ , let  $\alpha = \varepsilon^{1/(10C^{k+1}(k+1)!)}$ . Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with  $w_i \geq \alpha$  for all  $i \in [k]$ , and let  $\mu$  and  $\Sigma$  be the mean and covariance of  $\mathcal{M}$  such that  $r = \text{rank}(\Sigma)$  and for all  $i, j \in [k]$ ,  $\|\Sigma^{+1/2}(\Sigma_i - \Sigma_j)\Sigma^{+1/2}\|_F \leq 1/\sqrt{\alpha}$ . Let  $X$  be a set of points satisfying Condition 2.47 with respect to  $\mathcal{M}$  for some parameters  $(\gamma, t)$ . Given a set  $Y$ , an  $\varepsilon$ -corrupted version of  $X$ , of size  $n \geq n_0 = d^{O(1)}$ , there exists an algorithm that takes  $Y$  as input and in time  $n^{O(1)}$  outputs estimators  $\hat{\mu}$  and  $\hat{\Sigma}$  such that  $\hat{\Sigma} = \hat{U}\hat{\Lambda}\hat{U}^\top$  is the eigenvalue decomposition, where  $\hat{U} \in \mathbb{R}^{n \times r}$  has orthonormal columns and  $\hat{\Lambda} \in \mathbb{R}^{r \times r}$  is a diagonal matrix. Further, we can obtain  $n$  samples  $Y'$  by applying the affine transformation  $y_i \rightarrow \hat{U}^\top \hat{\Sigma}^{+1/2}(y_i - \hat{\mu})$  to each sample, such that a  $(1 - \varepsilon)$ -fraction have mean  $\mu'$  and covariance  $\Sigma'$  satisfying*

1.  $\|\mu'\|_2 \leq O\left(\left(1 + \frac{\sqrt{\varepsilon k}}{\alpha}\right) \sqrt{\varepsilon/\alpha}\right),$

2.  $\left(\frac{1}{1+(k\sqrt{\varepsilon}/\alpha)}\right) I_r \leq \Sigma' \leq \left(\frac{1}{1-(k\sqrt{\varepsilon}/\alpha)}\right) I_r,$
3.  $\|\Sigma' - I_r\|_F \leq O(\sqrt{\varepsilon}k/\alpha),$

where  $I_r$  is the  $r$ -dimensional Identity matrix, and the remaining points are arbitrary. Let  $X'$  be the set obtained by  $\hat{U}^\top \hat{\Sigma}^{+2} (x_i - \hat{\mu})$ . Then,  $X'$  satisfies Condition 2.47 with respect to  $\sum_{i=1}^k w_i \mathcal{N}\left(\hat{U}^\top \hat{\Sigma}^{+2}(\mu_i - \hat{\mu}), \hat{U}^\top \hat{\Sigma}^{+2} \Sigma_i \hat{\Sigma}^{+2} \hat{U}\right)$  and parameters  $(\gamma, t)$ , and  $Y'$  is an  $\varepsilon$ -corruption of  $X'$ .

*Proof.* For any  $t' \in \mathbb{N}$ , it follows from Corollary 2.35 that  $\mathcal{M}$  has  $2t'$ -certifiably  $(4/\alpha)$ -hypercontractive degree-2 polynomials, since  $w_i \geq \alpha$  for all  $i$ . Next, Lemma 6.6 implies that the uniform distribution over  $X$  also has  $2t'$ -certifiably  $(8/\alpha)$ -hypercontractive degree-2 polynomials and for  $d \times d$ -matrix-valued indeterminate  $Q$ ,

$$\left| \frac{Q}{2} \left\{ \mathbb{E}_{\mathcal{M}} \left( x^\top Q x - \mathbb{E}_{\mathcal{M}} x^\top Q x \right)^2 \leq O(1/\alpha) \|\Sigma^{1/2} Q \Sigma^{1/2}\|_F^2 \right\} \right|.$$

Then, it follows from Fact 2.38 that if  $\frac{16}{\alpha} t' \varepsilon^{1-4/t'} \ll 1$ , there exists an algorithm that runs in time  $n^{O(t')}$  and outputs estimates  $\hat{\mu}$  and  $\hat{\Sigma}$  satisfying:

1.  $\|\Sigma^{+2}(\mu - \hat{\mu})\|_2 \leq O(t'/\alpha)^{1/2} \varepsilon^{1-1/t'},$
2.  $(1 - (k/\alpha)\varepsilon^{1-2/t'}) \Sigma \leq \hat{\Sigma} \leq (1 + (k/\alpha)\varepsilon^{1-2/t'}) \Sigma$  and,
3.  $\|\Sigma^{+2}(\hat{\Sigma} - \Sigma)\|_F \leq (t'/\alpha) O(\varepsilon^{1-1/t'}).$

Setting  $t' = 2$ , compute  $\hat{\Sigma} = \hat{U} \hat{\Lambda} \hat{U}^\top$ , the eigendecomposition of  $\hat{\Sigma}$ , such that  $\hat{U} \in \mathbb{R}^{n \times r}$  has orthonormal columns, where  $r \leq d$  is the rank of  $\hat{\Sigma}$  and  $\hat{\Lambda} \in \mathbb{R}^{r \times r}$  is a diagonal matrix. Similarly, let  $\Sigma = U \Lambda U^\top$  be the eigendecomposition of  $\Sigma$ . We apply the affine transformation  $y_i \rightarrow \hat{U}^\top \hat{\Sigma}^{+2} (y_i - \hat{\mu})$  to each sample and thus we can assume throughout the rest of our argument that we have access to  $\varepsilon$ -corrupted samples from a  $k$ -mixture of Gaussians with mean  $\mu' = \hat{U}^\top \hat{\Sigma}^{+2}(\mu - \hat{\mu})$  and covariance  $\Sigma' = \hat{U}^\top \hat{\Sigma}^{+2} \Sigma \hat{\Sigma}^{+2} \hat{U}$ . Then, we have that

$$\begin{aligned} \|\mu'\|_2 &= \left\| \hat{U}^\top \hat{\Sigma}^{+2}(\mu - \hat{\mu}) \right\|_2 \leq \left\| \hat{U}^\top \right\|_{\text{op}} \left\| \hat{\Sigma}^{+2}(\mu - \hat{\mu}) \right\|_2 \\ &\leq O\left( \left( 1 + \frac{\sqrt{\varepsilon}k}{\alpha} \right) \sqrt{\varepsilon/\alpha} \right), \end{aligned}$$

where the last inequality follows from (1) and (2). It also follows from (2) that

$$\left( \frac{1}{1 + (k\sqrt{\varepsilon}/\alpha)} \right) \hat{\Sigma} \leq \Sigma \leq \left( \frac{1}{1 - (k\sqrt{\varepsilon}/\alpha)} \right) \hat{\Sigma}. \quad (6.1)$$

Multiplying out (6.1) with  $\hat{U}^\top \hat{\Sigma}^{+2}$  on the left and  $\hat{\Sigma}^{+2} \hat{U}$  on the right, we have

$$\left( \frac{1}{1 + (k\sqrt{\varepsilon}/\alpha)} \right) \hat{U}^\top \hat{\Sigma}^{+2} \hat{\Sigma} \hat{\Sigma}^{+2} \hat{U} \leq \Sigma' \leq \left( \frac{1}{1 - (k\sqrt{\varepsilon}/\alpha)} \right) \hat{U}^\top \hat{\Sigma}^{+2} \hat{\Sigma} \hat{\Sigma}^{+2} \hat{U}.$$

Observe that (2) implies that the rank of  $\hat{\Sigma}$  and  $\Sigma$  is the same, and thus  $\hat{U}^\top \hat{\Sigma}^{+1/2} \hat{\Sigma} \hat{\Sigma}^{+1/2} \hat{U} = I_r$ , where  $I_r$  is the  $r$ -dimensional Identity matrix. Finally, we have that

$$\begin{aligned} \|\Sigma' - I_r\|_F &= \left\| \hat{U}^\top \hat{\Sigma}^{+1/2} \Sigma \hat{\Sigma}^{+1/2} \hat{U} - \hat{U}^\top \hat{\Sigma}^{+1/2} \hat{\Sigma} \hat{\Sigma}^{+1/2} \hat{U} \right\|_F \leq \left\| \hat{U} \hat{\Lambda}^{-1/2} \hat{U}^\top (\Sigma - \hat{\Sigma}) \hat{U} \hat{\Lambda}^{-1/2} \hat{U}^\top \right\|_F \\ &= \left\| \hat{U} \hat{\Lambda}^{-1/2} \Lambda^{1/2} \Lambda^{-1/2} \hat{U}^\top (\Sigma - \hat{\Sigma}) \hat{U} \hat{\Lambda}^{-1/2} \Lambda^{1/2} \hat{\Lambda}^{-1/2} \hat{U}^\top \right\|_F \\ &\leq \left\| \hat{\Lambda}^{-1/2} \Lambda^{1/2} \right\|_{\text{op}}^2 \left\| \Sigma^{+1/2} (\hat{\Sigma} - \Sigma) \Sigma^{+1/2} \right\|_F \\ &\leq \mathcal{O}\left(\sqrt{\varepsilon} k / \alpha\right), \end{aligned}$$

where we use that  $\hat{\Lambda}^{-1/2} = \hat{\Lambda}^{-1/2} \Lambda^{1/2} \Lambda^{-1/2}$ , the sub-multiplicative property of the Frobenius norm, the column span  $U$  and  $\hat{U}$  is identical (see (2)), and the Frobenius recovery guarantee in (3).

Finally, it follows from Lemma 2.48 that Condition 2.47 is affine invariant and is thus preserved under  $x_i \rightarrow \hat{U}^\top \hat{\Sigma}^{-1/2} (x_i - \hat{\mu})$ , for  $i \in [n]$ , with parameters  $(\gamma, t)$ .  $\square$

The above robust isotropic transformation lemma allows us to obtain a covariance that is close to the identity matrix in a full-dimensional subspace (potentially smaller than the input dimension). Therefore, we will subsequently drop the subscript for the dimension, wherever it is clear from the context.

Next, we show that whenever the minimum mixing weight is sufficiently larger than the fraction of outliers, and a pair of components is covariance separated, we can partially cluster the samples.

**Lemma 6.8** (Non-negligible Weight and Covariance Separation). *Given  $0 < \varepsilon < 1/k^{k^{\mathcal{O}(k^2)}}$  and  $k \in \mathbb{N}$ , let  $\alpha = \varepsilon^{1/(10C^{k+1}(k+1)!)}$ . Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with mixture covariance  $\Sigma$  such that  $w_i \geq \alpha$  for all  $i \in [k]$  and there exist  $i, j \in [k]$  such that  $\left\| \Sigma^{+1/2} (\Sigma_i - \Sigma_j) \Sigma^{+1/2} \right\|_F > 1/\sqrt{\alpha}$ . Further, let  $X$  be a set of points satisfying Condition 2.47 with respect to  $\mathcal{M}$  for some parameters  $\gamma \leq \varepsilon d^{-8k} k^{-Ck}$ , for a sufficiently large constant  $C$ , and  $t \geq 8k$ . Let  $Y$  be an  $\varepsilon$ -corrupted version of  $X$  of size  $n \geq n_0 = (dk)^{\Omega(1)} / \varepsilon$ , Algorithm 4.4 partitions  $Y$  into  $Y_1, Y_2$  in time  $n^{\mathcal{O}(1)}$  such that with probability at least  $\alpha^{k \log(k/\alpha)}$  there is a non-trivial partition of  $[k]$  into  $Q_1 \cup Q_2$  so that letting  $\mathcal{M}_j$  be a distribution proportional to  $\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then  $Y_j$  is an  $\mathcal{O}\left(\varepsilon^{1/(10C^{k+1}(k+1)!)}\right)$ -corrupted version of  $\bigcup_{i \in Q_j} X_i$  satisfying Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\mathcal{O}(k\gamma/W_j), t)$ .*

*Proof.* We run Algorithm 4.4 with sample set  $Y$ , number of components  $k$ , the fraction of outliers  $\varepsilon$  and the accuracy parameter  $\beta$ . Since  $X$  satisfies Condition 2.47, we can set  $t' \geq 24$ ,  $\beta = \alpha^{t'/4-4} k^{t'} (t')^{2t'} \leq \alpha$  in Theorem 4.3. Then, by assumption, there exist  $i, j$  such that

$$\left\| \Sigma^{+1/2} (\Sigma_i - \Sigma_j) \Sigma^{+1/2} \right\|_F > \frac{1}{\sqrt{\alpha}} = \Omega\left(\frac{k^2 (t')^4}{(\beta \alpha^4)^{2/t'}}\right).$$

We observe that we also satisfy the other preconditions for Theorem 4.3, since  $n \geq (dk)^{\Omega(1)} / \varepsilon$ .

Then, Theorem 4.3 implies that with probability at least  $\alpha^{k \log(k/\alpha)}$ , the set  $Y$  is partitioned in two sets  $Y_1$  and  $Y_2$  such that there is a non-trivial partition of  $[k]$  into  $Q_1 \cup Q_2$  so that letting  $\mathcal{M}_j$

be a distribution proportional to  $\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then  $Y_j$  is an  $\mathcal{O}\left(\varepsilon^{1/(10C^{k+1}(k+1)!)}\right)$ -corrupted version of  $\bigcup_{i \in Q_j} X_i$ . By Lemma 2.50,  $\bigcup_{i \in Q_j} X_i$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\mathcal{O}(k\gamma/W_j), t)$ .  $\square$

When the mixture is not covariance separated and nearly isotropic, we can obtain a small list of hypotheses such that one of them is close to the true parameters, via tensor decomposition.

**Lemma 6.9** (Mixture is List-decodable). *Given  $0 < \varepsilon < 1/k^{k^{\mathcal{O}(k^2)}}$  let  $\alpha = \varepsilon^{1/(10C^{k+1}(k+1)!)}$ . Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with mixture mean  $\mu$  and mixture covariance  $\Sigma$ , such that  $\|\mu\|_2 \leq \mathcal{O}(\sqrt{\varepsilon/\alpha})$ ,  $\|\Sigma - I\|_F \leq \mathcal{O}(\sqrt{\varepsilon/\alpha})$ ,  $w_i \geq \alpha$  for all  $i \in [k]$ , and  $\|\Sigma_i - \Sigma_j\|_F \leq 1/\sqrt{\alpha}$  for any pair of components, and let  $X$  be a set of points satisfying Condition 2.47 with respect to  $\mathcal{M}$  for some parameters  $\gamma = \varepsilon d^{-8k} k^{-Ck}$ , for a sufficiently large constant  $C$ , and  $t = 8k$ . Let  $Y$  be an  $\varepsilon$ -corrupted version of  $X$  of size  $n$ , Algorithm 3.2 outputs a list  $L$  of hypotheses of size  $\exp(1/\varepsilon^2)$  in time  $\text{poly}(|L|, n)$  such that if we choose a hypothesis  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$  uniformly at random,  $\|\mu_i - \hat{\mu}_i\|_2 \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$  and  $\|\Sigma_i - \hat{\Sigma}_i\|_F \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$  for all  $i$  with probability at least  $\exp(-1/\varepsilon^2)$ .*

*Proof.* Recall we run Algorithm 3.2 on the samples  $Y$ , the number of clusters  $k$ , the fraction of outliers  $\varepsilon$  and the minimum weight  $\alpha = \varepsilon^{1/(10C^{k+1}(k+1)!)}$ . Next, we show that the preconditions of Theorem 3.1 are satisfied. First, the upper bounds on  $\|\mu\|_2$  and  $\|\Sigma - I\|_F$  imply  $\sum_{i \in [k]} w_i (\Sigma_i + \mu_i \mu_i^\top) = \Sigma + \mu \mu^\top \leq (1 + \mathcal{O}(\sqrt{\varepsilon/\alpha}))I$ . Since the LHS is a conic combination of PSD matrices, it follows that for all  $i \in [k]$ ,  $\mu_i \mu_i^\top \leq \frac{1}{\alpha} (1 + \mathcal{O}(\sqrt{\varepsilon/\alpha}))I$ , and thus  $\|\mu_i \mu_i^\top\|_F \leq \frac{2}{\alpha}$ . Next, we can write:

$$\begin{aligned} \|\Sigma_i - I\|_F &\leq \left\| \Sigma_i - (\Sigma + \mu \mu^\top) \right\|_F + \|\Sigma - I\|_F + \|\mu \mu^\top\|_F \\ &= \left\| \Sigma_i - \sum_{j \in [k]} w_j (\Sigma_j + \mu_j \mu_j^\top) \right\|_F + \frac{\sqrt{\varepsilon}k}{\alpha} + \frac{\varepsilon}{\alpha} \\ &\leq \left\| \sum_{j \in [k]} w_j (\Sigma_i - \Sigma_j) \right\|_F + \frac{2}{\alpha} + \frac{\sqrt{\varepsilon}k}{\alpha} + \frac{\varepsilon}{\alpha} \\ &\leq \frac{4}{\alpha}, \end{aligned}$$

where the first and the third inequalities follow from the triangle inequality and the upper bound on  $\|\mu_i \mu_i^\top\|_F$ , and the last inequality follows from the assumption that  $\|\Sigma_i - \Sigma_j\|_F \leq 1/\sqrt{\alpha}$  for every pair of covariances  $\Sigma_i, \Sigma_j$ . So, we can set  $\Delta = 4/\alpha$  in Theorem 3.1. Then, given the definition of  $\alpha$ , we have that

$$\eta = 2k^{4k} \mathcal{O}(1 + \Delta/\alpha)^{4k} \sqrt{\varepsilon} = \mathcal{O}\left(\varepsilon^{2/5}\right)$$

and  $1/\varepsilon^2 \geq \log(1/\eta)(k + 1/\alpha + \Delta)^{4k+5}/\eta^2$ . Therefore, Algorithm 3.2 outputs a list  $L$  of hypotheses such that  $|L| = \exp(1/\varepsilon^2)$ , and with probability at least 0.99,  $L$  contains a hypothesis that satisfies



the following: for all  $i \in [k]$ ,

$$\begin{aligned} \|\hat{\mu}_i - \mu_i\|_2 &= \mathcal{O}\left(\frac{\Delta^{1/2}}{\alpha}\right) \eta^{G(k)} = \mathcal{O}\left(\varepsilon^{-1/(20C^{k+1}(k+1)!)} \cdot \varepsilon^{1/(10C^{k+1}(k+1)!)}\right) = \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right) \text{ and} \\ \|\hat{\Sigma}_i - \Sigma_i\|_F &= \mathcal{O}(k^4) \frac{\Delta^{1/2}}{\alpha} \eta^{G(k)} = \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right). \end{aligned} \quad (6.2)$$

Then if we choose a hypothesis in  $L$  uniformly at random, the probability that we choose the hypothesis satisfying (6.2) is at least  $1/|L| = \exp(-1/\varepsilon^2)$ .  $\square$

Finally, if the mixture has a covariance matrix with small variance along any direction, we can further cluster the points by projecting the mixture along that direction.

**Lemma 6.10** (Spectral Separation of Thin Components). *Given  $0 < \varepsilon < 1/k^{k^{O(k^2)}}$ , let  $\alpha = \varepsilon^{1/(10C^{k+1}(k+1)!)}$ . Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with mixture covariance  $\Sigma$  such that  $\|\Sigma - I\|_F \leq \mathcal{O}(\sqrt{\varepsilon}k/\alpha)$ , and let  $X$  be a set of points satisfying Condition 2.47 with respect to  $\mathcal{M}$  for some parameters  $(\gamma, t)$ . Given a set  $Y$  being an  $\varepsilon$ -corrupted version of  $X$  of size  $n$ , and estimates  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$ , such that  $\|\mu_i - \hat{\mu}_i\|_2 \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$ ,  $\|\Sigma_i - \hat{\Sigma}_i\|_F \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$ , suppose there exists a unit vector  $v \in \mathbb{R}^d$  such that  $v^\top \hat{\Sigma}_s v \leq \mathcal{O}\left(\varepsilon^{1/(40C^{k+1}(k+1)!)}\right)$ , for some  $s \in [k]$ . Then, there is an algorithm that efficiently partitions  $Y$  into  $Y_1$  and  $Y_2$  such that there is a non-trivial partition of  $[k]$  into  $Q_1 \cup Q_2$  so that letting  $\mathcal{M}_j$  be a distribution proportional to  $\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then  $Y_j$  is an  $\left(\mathcal{O}(k^2\gamma) + \mathcal{O}\left(\varepsilon^{1/(80C^{k+1}(k+1)!)} / W_j\right)\right)$ -corrupted version of  $\bigcup_{i \in Q_j} X_i$  satisfying Condition 2.47 with respect to  $\mathcal{M}_j$  with parameter  $(\mathcal{O}(k\gamma/W_j), t)$ .*

*Proof.* We run the algorithm from Lemma 5.1 with the input being the samples  $Y$ , the current hypothesis  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$ , and the minimum eigenvalue  $\eta = \mathcal{O}\left(\varepsilon^{1/(40C^{k+1}(k+1)!)}\right)$ . Observe that the mixture covariance satisfies  $\Sigma \geq (1 - \mathcal{O}(\sqrt{\varepsilon}k/\alpha))I \geq I/2$  and the upper bound on means and covariance is  $\delta = \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}n\right)$  by assumption. Therefore, we satisfy the preconditions of Lemma 5.1. Thus, we obtain a partition  $Y_1, Y_2$  such that there is a non-trivial partition of  $[k]$  into  $Q_1 \cup Q_2$  so that letting  $\mathcal{M}_j$  be a distribution proportional to  $\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then it follows from Lemma 2.50 that  $Y_j$  is an  $\left(\mathcal{O}(k^2\gamma) + \mathcal{O}\left(\varepsilon^{1/(80C^{k+1}(k+1)!)} / W_j\right)\right)$ -corrupted version of  $\bigcup_{i \in Q_j} X_i$  satisfying Condition 2.47 with respect to  $\mathcal{M}_j$  with parameter  $(\mathcal{O}(k\gamma/W_j), t)$ .  $\square$

## 6.2 Proof of the Main Theorem

We are now ready to complete the proof of Theorem 6.1.

*Proof of Theorem 6.1.* We divide the proof into two parts: first we show that Algorithm 6.4 outputs a hypothesis  $\widehat{\mathcal{M}} = \sum_{i \in [k]} \hat{w}_i \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i)$  such that  $\widehat{\mathcal{M}}$  and  $\mathcal{M}$  are  $\mathcal{O}(\varepsilon^{c_k})$ -close in total variation distance with probability at least  $\exp(-\mathcal{O}(k)/\varepsilon^2)$ ; then we show that Algorithm 6.3 outputs a  $k$ -mixture of Gaussians  $\widehat{\mathcal{M}}$  such that  $\widehat{\mathcal{M}}$  and  $\mathcal{M}$  are  $\mathcal{O}_k(\varepsilon^{c_k})$ -close in total variation distance with probability 0.99.



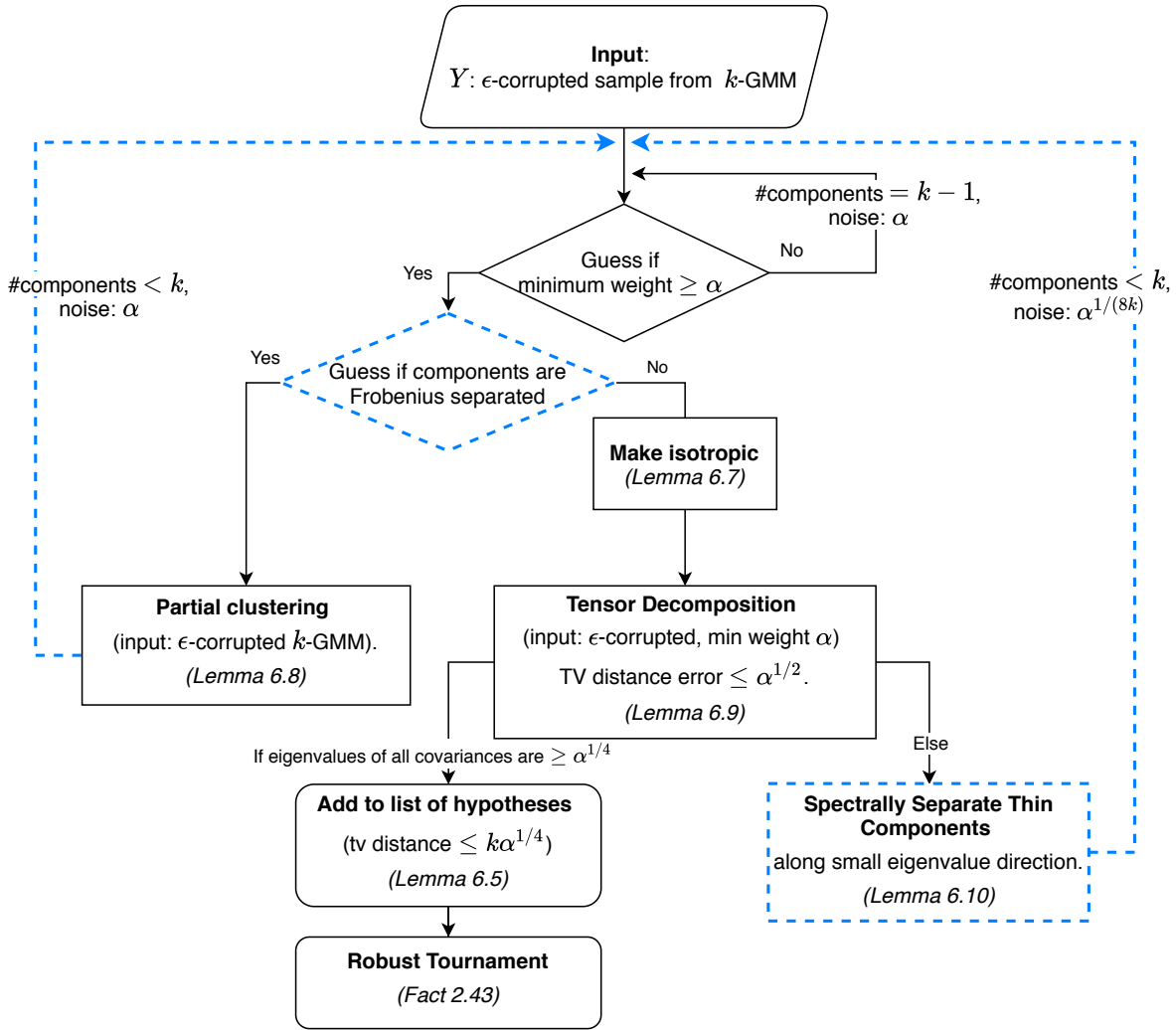


Figure 1: If we assume a  $1/\text{poly}(k)$  lower bound on minimum weight, then we can skip all blue steps above; the partial clustering is carried out till it can no longer be done within a cluster and then followed by the tensor decomposition step.

We proceed the first part by induction on  $k$ . Let  $c_k = \frac{1}{(100)^k C^{(k+1)! \text{sf}(k+1)k!}}$  be a scalar that only depends on  $k$ , where  $C > 0$  is a sufficiently large universal constant.

**Induction Hypothesis:** Let  $X$  be a set of points satisfying Condition 2.47 with respect to a  $k$ -mixture of Gaussians  $\mathcal{M}$  for some parameters  $\gamma = \varepsilon d^{-8k} k^{-C'k}$ , where  $C'$  is a sufficiently large constant and  $t = 8k + 48$ . Given a set  $Y$  being an  $\varepsilon$ -corrupted version of  $X$  of size  $n$ , the outlier parameter  $\varepsilon$  and the component-number parameter  $k$ , Algorithm 6.3 returns a  $k$ -mixture of Gaussians  $\widehat{\mathcal{M}}$  such that  $\widehat{\mathcal{M}}$  and  $\mathcal{M}$  are  $\mathcal{O}_k(\varepsilon^{c_k})$ -close in total variation distance with probability  $\exp(-(3k - 2)/\varepsilon^2)$ .

**Base Case:** For  $k = 1$ , the algorithm returns the single Gaussian with mean  $\hat{\mu}$  and  $\hat{\Sigma}$  at Step 2.

Suppose the true Gaussian is  $N(\mu, \Sigma)$ . It follows from the proof of Lemma 6.7,

$$\|\Sigma^{+/2} (\hat{\mu} - \mu)\|_2 = \|\Sigma^{+/2} (\hat{\mu} - \mu)\|_2 \leq \mathcal{O}(\sqrt{\varepsilon})$$

and

$$\left\| \Sigma^{+/2} (\hat{\Sigma} - \Sigma) \Sigma^{+/2} \right\|_F \leq \mathcal{O}(\sqrt{\varepsilon}),$$

and thus it follows from Fact 2.1 that the total variation distance between the hypothesis Gaussian and the true Gaussian is at most  $\mathcal{O}(\sqrt{\varepsilon})$ . We can then conclude that the base case is true.

**Inductive Step:** We assume that our induction hypothesis holds for any  $m < k$  and then prove that the induction hypothesis holds for  $k$ .

**Small Clusters Can be Treated as Noise.** Conditioning on the base case being true, we begin by guessing whether the minimum weight is less than  $\varepsilon^{1/(10C^{k+1}(k+1)!)}$  with equal probability.

Let  $w_{\min} = \min_i w_i$ . If  $w_{\min} \leq \varepsilon^{1/(10C^{k+1}(k+1)!)}$ , our algorithm takes step 1 with probability 0.5. In this case, we treat the smallest component as noise and recurse on the set of samples  $Y$ . We set the number of components to be  $k - 1$  and the fraction of outliers being  $\varepsilon + \varepsilon^{1/(10C^{k+1}(k+1)!)} \leq 2\varepsilon^{1/(10C^{k+1}(k+1)!)}$ . By Lemma 2.50,  $Y$  is an  $2\varepsilon^{1/(10C^{k+1}(k+1)!)}$ -corrupted version of a set satisfying Condition 2.47 with respect to a  $(k - 1)$ -mixture for parameters  $\gamma = \mathcal{O}(k\varepsilon d^{-8k} k^{-C'k} / (1 - w_{\min})) \leq \varepsilon d^{-8(k-1)} (k-1)^{-C'(k-1)}$  and  $t = 8k + 48$ . Thus applying the inductive hypothesis to  $Y$ , we learn the mixture up to total variation distance  $\mathcal{O}_k\left(\left(2\varepsilon^{1/(10C^{k+1}(k+1)!)}\right)^{C^{k-1}}\right) \leq \mathcal{O}_k(\varepsilon^{C^k})$  with probability  $0.5 \exp(-(3(k-1) - 2)/\varepsilon^2) \geq \exp(-(3k - 2)/\varepsilon^2)$ . Now we may assume for all  $i \in [k]$ ,  $w_i \geq \varepsilon^{1/(10C^{k+1}(k+1)!)}$ .

**Mixture is Covariance Separated.** Let  $\alpha = \varepsilon^{1/(10C^{k+1}(k+1)!)}$  and  $\psi_1 = \{\exists \mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j) \mid \|\Sigma_i - \Sigma_j\|_F > \alpha^{-1/2}\}$  be the event that the samples were drawn from a mixture that is *covariance separated*. First, consider the case where  $\psi_1$  is true. We will run 3(a) with probability 0.5. Then it follows from Lemma 6.8 that  $Y$  can be partitioned into  $Y_1$  and  $Y_2$  in time  $d^{O(1)}$ , such that they both have at least one component and the fraction of outliers in each set  $Y_1, Y_2$  is at most  $\varepsilon^{1/(10C^{k+1}(k+1)!)}$  with probability  $\alpha^{O(k \log(k/\alpha))}$ . Then, we can guess the number of components in  $Y_1$  and we will be correct with probability  $1/k$ . Conditioned on our guess being correct, let  $Y_1$  consist of  $k_1$  components and  $Y_2$  consist of  $k_2$  components and  $k_1 + k_2 = k$ .

Let  $Q_1 \cup Q_2$  be the non-trivial partition of  $[k]$  in Lemma 6.8,  $\mathcal{M}_j$  be a distribution proportional to  $\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then By Lemma 2.50,  $Y_j$  is an  $\mathcal{O}\left(\varepsilon^{1/(10C^{k+1}(k+1)!)}\right)$ -corrupted version of  $\bigcup_{i \in Q_j} X_i$  satisfying Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $\gamma = \mathcal{O}(k\varepsilon d^{-8k} k^{-C'k} / \alpha) \leq \varepsilon d^{-8k_j} (k_j)^{-C'k_j}$  and  $t = 8k + 48$ . Then, applying the inductive hypothesis on  $Y_j$  for  $j = 1, 2$ , with number of components  $k_j$ , we can learn the mixtures  $\mathcal{M}_j$  up to total variation distance error  $\mathcal{O}_k\left(\varepsilon^{C_{k_j}/(10C^{k+1}(k+1)!)}\right)$  with probability  $\exp(-(3k_j - 2)/\varepsilon^2)$ . Finally if this is the case, we combine the two hypotheses on  $Y_1, Y_2$  by multiplying each weight in the hypothesis of  $Y_j$  by  $|Y_j|/|Y|$  and then taking union of two hypotheses. Then our combining method gives a final output that learns our full hypothesis to total variation distance error  $\mathcal{O}_k\left(\varepsilon^{C_{k_1}/(10C^{k+1}(k+1)!)}\right) + \mathcal{O}_k\left(\varepsilon^{C_{k_2}/(10C^{k+1}(k+1)!)}\right) \leq$

$\mathcal{O}_k(\varepsilon^{c_k})$  with probability at least  $0.5 \cdot 0.5 \cdot \frac{1}{k} \cdot \alpha^{O(k \log(k/\alpha))} \exp(-(3k_1 - 2)/\varepsilon^2) \exp(-(3k_2 - 2)/\varepsilon^2) \geq \exp(-(3k - 2)/\varepsilon^2)$ .

**Mixture is not Covariance Separated.** Next, consider the case where  $\psi_1$  is false. With probability 0.5, the algorithm guesses correctly and executes Step 2. Since the mixture is not covariance separated, we satisfy the preconditions of Lemma 6.7, and after applying the transformation in Step 2,  $\Sigma$ , the covariance of the mixture  $\mathcal{M}$ , is  $\sqrt{\varepsilon}k/\alpha$ -close to the  $r$ -dimensional identity, where  $r$  is the rank of  $\Sigma$ . However, since we obtain the subspace exactly, we can simply project all samples on the subspace and we drop the  $r$  in the subsequent exposition.

Let  $X'$  be the set of points obtained by applying the Affine transformation from Step 2 as defined in Lemma 6.7. Then,  $X'$  satisfies Condition 2.47 with respect to a nearly isotropic mixture and parameters  $\gamma = \varepsilon d^{-8k} k^{-C^k}$  and  $t = 8k + 48$  so that we can continue the algorithm with  $X'$ . Whenever we return a hypothesis in the following steps, we will first apply the inverse of the transformation on our estimates  $\hat{\mu}_i$  and  $\hat{\Sigma}_i$ . Since total variation distance is affine invariant, we have the same error guarantee in total variation distance after applying the transformation. From now on, we reduce to the case where  $\Sigma$  is  $\sqrt{\varepsilon}k/\alpha$ -close to the Identity.

There is a 50% chance our algorithm runs Step 3(b) and we will analyze the remainder of this case under that assumption. It follows from Lemma 6.9 that we obtain a hypothesis  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$  such that  $\|\mu_i - \hat{\mu}_i\|_2 \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$  and  $\|\Sigma_i - \hat{\Sigma}_i\|_F \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$  with probability  $\exp(-1/\varepsilon^2)$ . Conditioned on the hypothesis being correct, we now split into two cases: either all eigenvalues of all the estimated covariances are large (in which case we obtain total variation distance guarantees), or there is a direction along which we can project and cluster further.

**Covariance Estimates have Large Eigenvalues.** For the hypothesis  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$  from the last step, we compute all the eigenvalues of the estimated covariance matrices,  $\hat{\Sigma}_i$ , for all  $i \in [k]$ . If, for all  $i \in [k]$ ,  $\lambda_{\min}(\hat{\Sigma}_i) \geq c\varepsilon^{1/(40C^{k+1}(k+1)!)}$ , we land in Step 3(b).i that we guess the mixing weights  $\hat{w}_i$  uniformly in the range  $[0, 1]$  and then we output the corresponding hypothesis  $\{\hat{w}_i, \hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$ . With probability at least  $\varepsilon^k$ ,  $\hat{w}_i$  are within  $\varepsilon$  of the true mixing weights. Under this condition, by Lemma 6.5, the mixture  $\widehat{\mathcal{M}} = \sum_{i \in [k]} \hat{w}_i \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i)$  is  $\mathcal{O}_k\left(\varepsilon^{1/(40C^{k+1}(k+1)!)}\right) \leq \mathcal{O}_k(\varepsilon^{c_k})$ -close to  $\mathcal{M}$  in total variation distance with probability  $0.5 \cdot 0.5 \cdot \varepsilon^k \cdot \exp(-1/\varepsilon^2) \geq \exp(-(3k - 2)/\varepsilon^2)$ .

**One Covariance Has a Small Eigenvalue.** Consider the case (Step 3(b).ii) where there exists a unit-norm direction  $v$  and an estimate  $\hat{\Sigma}_i$  such that  $v^\top \hat{\Sigma}_i v \leq c\varepsilon^{1/(40C^{k+1}(k+1)!)}$ . It then follows from Lemma 6.10 that we can partition  $Y$  into  $Y_1$  and  $Y_2$  such that each has at least one cluster and the total number of outliers in both  $Y_1$  and  $Y_2$  is at most  $\mathcal{O}\left(\varepsilon^{1/(80kC^{k+1}(k+1)!)}\right)n$ . If  $Y_1$  or  $Y_2$  has size less than  $\varepsilon^{1/(400kC^{k+1}(k+1)!)}n$ , then we can treat it as noise and get an additive  $\mathcal{O}\left(\varepsilon^{1/(400kC^{k+1}(k+1)!)}\right)$ -error in total variation distance. Otherwise, the fraction of outliers in both sets is at most  $\mathcal{O}\left(\left(\varepsilon^{1/(80kC^{k+1}(k+1)!)}n\right)/\left(\varepsilon^{1/(400kC^{k+1}(k+1)!)}n\right)\right) = \mathcal{O}\left(\varepsilon^{1/(100kC^{k+1}(k+1)!)}\right)$ . We then guess the number of components,  $k_1$ , in  $Y_1$  with success probability  $1/k$ . Let  $k_2 = k - k_1$  be the number of components

in  $Y_2$ . Then, conditioned on this event holding,  $Y_j$  is an  $\mathcal{O}\left(\varepsilon^{1/(100kC^{k+1}(k+1)!)}\right)$ -corrupted version of a set satisfying Condition 2.47 with respect to a mixture of  $k_j$  components with parameter  $\gamma = k\varepsilon d^{-8k}k^{-C^k}/\alpha \leq \varepsilon d^{-8(k_j)}(k_j)^{-C^k}$  and  $t = 8k + 48$ . We can apply the inductive hypothesis to  $Y_1$  with number of components  $k_1$  and fraction of outliers  $\mathcal{O}\left(\varepsilon^{1/(100kC^{k+1}(k+1)!)}\right)$ , and conclude that we learn the components of  $Y_1$  to total variation distance  $\mathcal{O}_k\left(\varepsilon^{c_{k_1}/(100kC^{k+1}(k+1)!)}\right)$  with probability  $\exp(-(3k_1 - 2)/\varepsilon^2)$ . A similar argument holds for  $Y_2$ . Finally if this is the case, we combine the two hypotheses on  $Y_1, Y_2$  by multiplying each weight by  $|Y_j|/|Y|$  and then taking union of two hypotheses. Then our combining method gives a final output that learns our full hypothesis to total variation distance error  $\mathcal{O}_k\left(\varepsilon^{c_{k_1}/(100kC^{k+1}(k+1)!)}\right) + \mathcal{O}_k\left(\varepsilon^{c_{k_2}/(100kC^{k+1}(k+1)!)}\right) + \mathcal{O}\left(\varepsilon^{1/(400kC^{k+1}(k+1)!)}\right) \leq \mathcal{O}_k(\varepsilon^{c_k})$  with probability at least  $0.5 \cdot 0.5 \cdot \frac{1}{k} \cdot \exp(-1/\varepsilon^2 - (3k_1 - 2)/\varepsilon^2 - (3k_2 - 2)/\varepsilon^2) \geq \exp(-(3k - 2)/\varepsilon^2)$ .

**Sample Size and Running Time of Algorithm 6.4** By Lemma 2.51, we need  $n \geq kt^{C^t}d^t/\gamma^3$  samples to generate  $X$  satisfying Condition 2.47 with parameters  $(\gamma, t)$ . We set  $\gamma = \varepsilon d^{-8k}k^{-C^k}$  and  $t = 8k + 48$ . Then  $n \geq n_0 = (8k)^{\mathcal{O}(k)}d^{\mathcal{O}(k)}/\varepsilon^3$ . The running time in each sub-routine we invoke is dominated by the running time of the tensor decomposition algorithm, and by Lemma 6.9 in the worst case this is  $\text{poly}(|L|, n) = \text{poly}(\exp(1/\varepsilon^2), d^{\mathcal{O}(k)}/\varepsilon^3) = d^{\mathcal{O}(k)}\exp(1/\varepsilon^2)$ .

This completes the first part of the proof.

**Aggregating Hypotheses.** We run Algorithm 6.4 repeatedly on set  $Y$  and add the return hypothesis into a list  $\mathcal{L}$  until with probability 0.99, there exists a hypothesis  $\widehat{M} \in \mathcal{L}$  such that  $\widehat{M}$  and  $M$  are  $\mathcal{O}_k(\varepsilon^{c_k})$ -close in total variation distance. Since Algorithm 6.4 outputs a correct mixture with probability  $\exp(-(3k - 2)/\varepsilon^2)$ , we will run Algorithm 6.4 for  $\exp(\mathcal{O}(k)/\varepsilon^2)$  times. Then the total running time is  $\exp(\mathcal{O}(k)/\varepsilon^2) \cdot d^{\mathcal{O}(k)}\exp(1/\varepsilon^2) = d^{\mathcal{O}(k)}\exp(\mathcal{O}(k)/\varepsilon^2)$ .

**Robust Tournament.** Then we need to run a robust tournament in order to find a hypothesis that is close to the true mixture in total variation distance. Fact 2.52 shows that we can do this efficiently only with access to an  $\varepsilon$ -corrupted set of samples of size  $\mathcal{O}_k(\log(1/\varepsilon)/\varepsilon^{2c_k})$ .

This completes the proof.  $\square$

## 7 More Efficient Robust Partial Cluster Recovery

In this section, we prove the following upgraded partial clustering theorem. In contrast to Theorem 4.3, here we obtain a probability of success that is inverse exponential in  $k$  instead of  $1/\alpha$ .

**Theorem 7.1** (Robust Partial Clustering in Relative Frobenius Distance). *Let  $0 \leq \varepsilon < \alpha/k \leq 1$  and  $t \in \mathbb{N}$ . There is an algorithm with the following guarantees: Let  $Y$  be an  $\varepsilon$ -corruption of a sample  $X$  of size  $n \geq (dk)^{C^t}/\varepsilon$  for a large enough constant  $C > 0$ , from  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  that satisfies Condition 2.47 with parameters  $2t$  and  $\gamma \leq \varepsilon d^{-8t}k^{-C^k}$ , for a large enough constant  $C > 0$ . Suppose further that  $w_i \geq \alpha > 2\varepsilon$  for each  $i \in [k]$ , and that for some  $t \in \mathbb{N}$ ,  $\beta > 0$  there exist  $i, j \leq k$  such that  $\|\Sigma^{+/2}(\Sigma_i - \Sigma_j)\Sigma^{+/2}\|_F^2 = \Omega((k^2t^4)/\beta^{2/t}\alpha^4)$ , where  $\Sigma$  is the covariance of the mixture  $\mathcal{M}$ . Then, for any*

$\eta \gg \sqrt{\varepsilon/\alpha}$ , the algorithm runs in time  $n^{O(t)}$ , and with probability at least  $2^{-O(k)}(1 - O(\eta/\alpha - \sqrt{\eta}))$  over the random choices of the algorithm, outputs a partition  $Y = Y_1 \cup Y_2$  satisfying:

1. **Partition respects clustering:** for each  $i$ ,  $\max \left\{ \frac{1}{w_{i1}} |Y_1 \cap X_i|, \frac{1}{w_{i2}} |Y_2 \cap X_i| \right\} \geq 1 - O(\sqrt{\eta}) - O\left(\frac{\beta}{\eta\alpha^2}\right)$ , where  $X_i \subset X$  corresponding to the points drawn from  $\mathcal{N}(\mu_i, \Sigma_i)$ .
2. **Partition is non-trivial:**  $\max_i \frac{1}{w_{i1}} |X_i \cap Y_1|, \max_i \frac{1}{w_{i2}} |X_i \cap Y_2| \geq 1 - O(\sqrt{\eta}) - O\left(\frac{\beta}{\eta\alpha^2}\right)$ .

## 7.1 Algorithm

Our algorithm will solve SoS relaxations of a polynomial inequality system. The indeterminates in this system are  $X'$  (that is intended to be the guess for the original uncorrupted sample), a cluster of size  $\alpha n$  within  $X'$  (indicated by  $z_i$ s) with mean  $\hat{\mu}$  and covariance matrix  $\hat{\Sigma}$  and  $\Pi$  (intended to be the square root of  $\hat{\Sigma}$ ). The input corrupted sample  $Y$  is a constant in this inequality system. Let  $U \in \mathbb{R}^{d \times d}$  and  $m, z \in \mathbb{R}^d$  also be indeterminates of the proof system. The system can be thought of as encoding the task of finding clusters  $\hat{C}$  within  $Y$  that satisfies certifiable hypercontractivity of degree 2 polynomials.

We present the constraints grouped together into meaningful categories below: The first set of constraints enforce that  $\hat{\Sigma}$  is the square of  $\Pi$ .

$$\text{Covariance Constraints: } \mathcal{A}_1 = \left\{ \begin{array}{l} \Pi = UU^\top \\ \Pi^2 = \hat{\Sigma} \end{array} \right\} \quad (7.1)$$

The intersection constraints force that  $X'$  intersects  $Y$  in all but an  $\varepsilon n$  points (and thus,  $2\varepsilon$ -close to unknown sample  $X$ ).

$$\text{Intersection Constraints: } \mathcal{A}_2 = \left\{ \begin{array}{l} \forall i \in [n], \quad m_i^2 = m_i \\ \quad \quad \quad \sum_{i \in [n]} m_i = (1 - \varepsilon)n \\ \forall i \in [n], \quad z_i(\tilde{y}_i - x'_i) = 0 \end{array} \right\} \quad (7.2)$$

The subset constraints enforce that  $z$  indicate a subset of size  $\alpha n$  of  $X'$ .

$$\text{Subset Constraints: } \mathcal{A}_3 = \left\{ \begin{array}{l} \forall i \in [n], \quad z_i^2 = z_i \\ \quad \quad \quad \sum_{i \in [n]} z_i = \alpha n \end{array} \right\} \quad (7.3)$$

Parameter constraints create indeterminates to stand for the covariance  $\hat{\Sigma}$  and mean  $\hat{\mu}$  of  $\hat{C}$  (indicated by  $z$ ).

$$\text{Parameter Constraints: } \mathcal{A}_4 = \left\{ \begin{array}{l} \frac{1}{\alpha n} \sum_{i=1}^n z_i (x'_i - \hat{\mu}) (x'_i - \hat{\mu})^\top = \hat{\Sigma} \\ \quad \quad \quad \frac{1}{\alpha n} \sum_{i=1}^n z_i x'_i = \hat{\mu} \end{array} \right\} \quad (7.4)$$

Certifiable Hypercontractivity :  $\mathcal{A}_4 =$

$$\left\{ \begin{array}{l} \forall t \leq 2s \quad \mathbb{E}_z(Q - \mathbb{E}_z Q)^{2t} \leq (Ct/\alpha)^t 2^{2t} \left( \mathbb{E}_z(Q - \mathbb{E}_z Q)^2 \right)^t \\ \mathbb{E}_z(Q - \mathbb{E}_z Q)^2 \leq 10 \left( \frac{1}{\alpha} \right)^2 \|Q\|_F^2 \end{array} \right\} \quad (7.5)$$

where we write  $\mathbb{E}_z Q$  as a shorthand for the polynomial  $\frac{1}{\alpha^n} \sum_i z_i Q(x_i)$  and  $\mathbb{E}_z(Q - \mathbb{E}_{X_r} Q)^{2j}$  for the polynomial  $\frac{1}{\alpha^n} \sum_i z_i (Q(x'_i) - \frac{1}{\alpha^n} \sum_{i \leq n} z_i Q(x'_i))^{2j}$  for any  $j$ . Note that  $Q$  is a  $d \times d$ -matrix valued indeterminate. Observe that  $Q$  itself can be eliminated from the system as is standard in several applications [KS17b, KS17a, HL18, BK20b] of SoS proofs in obtaining a succinct set of polynomial constraints (see Section 4.3 on ‘‘Succinct Representation of Constraints’’ in [FKP19] for an exposition).

**Algorithm 7.2** (Polynomial Time Partial Clustering).

**Given:** A sample  $Y$  of size  $n$ . An outlier parameter  $\varepsilon > 0$  and an accuracy parameter  $\eta > 0$ .

**Output:** A partition of  $Y$  into partial clustering  $Y_1 \cup Y_2$ .

**Operation:**

1. **Mean and Covariance Estimation:** Apply Robust Mean and Covariance Estimation (Fact 2.37) to estimate  $\hat{\mu}$  and  $\tilde{\Sigma}$  such that  $\frac{1}{2}\Sigma \leq \tilde{\Sigma} \leq 1.5\Sigma$  where  $\Sigma$  is the covariance of the uncorrupted input mixture.
2. **Approximate Isotropic Transformation:** For each  $y_i \in Y$ , let  $\tilde{y}_i = \tilde{\Sigma}^{+1/2}(y_i - \hat{\mu})$ . Let  $\tilde{Y} = \cup_{i \leq n} \tilde{y}_i$ .
3. **SDP Solving:** Find a pseudo-distribution  $\tilde{\zeta}$  satisfying  $\cup_{i=1}^4 \mathcal{A}_i$  such that  $\tilde{\mathbb{E}}_{\tilde{\zeta}} z_i \in \alpha \pm o_d(1)$  for every  $i$ . If no such pseudo-distribution exists, output fail.
4. **Rounding:** Let  $M = \tilde{\mathbb{E}}_{z \sim \tilde{\zeta}}[zz^\top]$ .
  - (a) **Generate candidate clusters:** For  $\ell = O(1/\alpha \log \eta/\alpha)$  times, draw a uniformly random  $i \in [n]$  and let  $\hat{C}_i = \{j \mid M(i, j) \geq \alpha^2/2\}$ . Let  $\mathcal{L} = \cup_{i \leq \ell} \hat{C}_i$ .
  - (b) **Candidate 2nd Moment Estimation:** For each  $\hat{C}_i \in \mathcal{L}$ , let  $S_i$  be the output of running robust 2nd moment estimation with Frobenius error (Lemma 7.7) on  $\hat{C}_i$  with outlier parameter  $\eta'_i = O(\frac{\varepsilon}{\alpha} + \frac{\beta}{\alpha^2 \eta})$ .
  - (c) **Merge candidate clusters:** For each  $i \leq \ell$ , find  $\mathcal{L}_i$  to be the collection of all  $j$  such that  $\|S_i - S_j\|_F \leq 2C\tau$  for a large enough constant  $C > 0$ . Set  $\hat{C}_i \cup \mathcal{L}_i = \hat{B}_i$ . Repeat on  $\mathcal{L} \setminus \{\hat{C}_i \cup i\}$ .
  - (d) **Output a union of a random subset of candidates:** For  $\mathcal{L}' = \cup_i \hat{B}_i$ , choose a uniformly random subset  $S$  of  $\mathcal{L}'$ , set  $Y_1 = \cup_{j \in S} \hat{B}_j$  and set  $Y_2 = Y \setminus Y_1$ .

**Analysis of Algorithm**

**Lemma 7.3** (Success of Step 1). *Let  $\tilde{\Sigma}$  be the output of the robust covariance estimation algorithm (Fact 2.37) applied to the input sample  $Y$  with outlier parameter  $\varepsilon$ . If  $Y$  is an  $\varepsilon$ -corruption of a sample  $X$  from a GMM with minimum weight  $\geq \alpha \geq \Omega(\sqrt{\varepsilon})$ , mixture mean  $\mu$  and covariance  $\Sigma$  satisfying Condition 2.47, then,*

$$0.5\Sigma \leq \tilde{\Sigma} \leq 1.5\Sigma, \\ \|\tilde{\Sigma}^{-1/2}(\mu - \hat{\mu})\|_2 \leq O(\sqrt{\varepsilon}/\alpha).$$

*Proof.* The lemma immediately follows by noting that GMMs with minimum weight  $\alpha$  are 4-certifiably  $1/\alpha$ -subgaussian (Fact 2.27) and  $\alpha \geq \Omega(\sqrt{\varepsilon})$ . □

**Lemma 7.4** (Simultaneous Intersection Bounds for Frobenius Separated Case). *Let  $X = X_1 \cup X_2 \cup \dots \cup X_k$  be a sample of size  $n \geq (dk)^{Ct}/\varepsilon$  for a large enough constant  $C > 0$ , from  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  that satisfies Condition 2.47 with parameters  $2t$  and  $\gamma \leq \varepsilon d^{-8t} k^{-Ck}$ . Suppose further that  $\|\mu_i\|_2 \leq \frac{2}{\alpha}$  for every  $i$ ,  $\|\Sigma_i\|_2 \leq \frac{1}{\alpha}$  for every  $i$  and the mixture mean  $\mu$ , covariance  $\Sigma$  satisfy  $\|\mu\|_2 \leq 1$  and  $0.5I \leq \Sigma \leq 1.5I$ . Let  $\tau = 10^8 \frac{C^6 \varepsilon^4}{\beta^{2t} \alpha^2}$ , for any  $\beta > 0$ . Then, given any  $\varepsilon$ -corruption  $Y$  of  $X$ , for every  $i, j$  such that  $\|\Sigma_i - \Sigma_j\|_F^2 \geq \Omega(\tau)$ ,*

$$\left\{ \bigcup_{i=1}^4 \mathcal{A}_i \right\} \Big|_{\frac{z}{2t}} \left\{ z'(X_i)z'(X_j) \leq \beta \right\},$$

where  $z'(X_i) = \frac{1}{w_i n} \sum_{j \in X_i} z_j \mathbf{1}(x_j = y_j)$  for every  $i$ .

*Proof of Theorem 7.1.* First, since  $Y$  is an  $\varepsilon$ -corruption of a sample  $X$  from a GMM such that  $X$  satisfies Condition 2.47, our robust mean and covariance estimation procedure (Step 1) applied to the mixture succeeds and recovers an estimate of the covariance that is multiplicative  $1 \pm 0.5$ -factor approximation in Löwner order. Thus, for the rest of the analysis, we can assume that the smallest and largest eigenvalue of the mixture covariance are at least 0.5 and at most 1.5. Since each component has weight at least  $\alpha$ , this means that each of the constitute component covariance can now be assumed to have a spectral norm at most  $1.5/\alpha$ .

Next, by an argument similar to the one presented in the proof of Theorem 4.2, the convex program we wrote is approximately solvable in polynomial time and is feasible whenever the uncorrupted sample  $X$  satisfies Condition 2.47. The only change here is in the certifiable hypercontractivity constraints where instead of the RHS of the bounded variance constraint is stated in terms of  $\|Q\|_F^2$  instead of  $\|\Pi Q \Pi\|_F^2$  with an additional slack of  $O(1/\alpha^2)$ . This modified constraint is satisfied by all true clusters by an application of Lemma 2.25 since each of their covariance has spectral norm at most  $1.5/\alpha$ .

**Rounding** Let  $M = \tilde{\mathbb{E}}_{\xi} z z^\top$ . Then, by an argument similar to the proof of Theorem 4.2, we can conclude:

1.  $o_d(1) + \alpha \geq M(i, j) \geq 0$ .
2.  $\sum_{j=1}^n M(i, j) \geq (\alpha^2 - o_d(1))n$  for every  $i$ .
3. For every  $i$ , let  $B_i$  be the set of “large entries”: i.e.  $j$  such that  $M(i, j) \geq \alpha^2/2$ . Then,  $|B_i| \geq \alpha n/2$ .



In the following, let  $M_i$  denote the  $i$ -th row of  $M$  and  $\|M_i\|_1$  for the sum of the non-negative entries of the vector  $M_i$ .

**Candidate Clusters** For every  $i$ , let  $F_i \subseteq [k]$  be the set of all  $i' \in [k]$  such that  $\|\Sigma_i - \Sigma_{i'}\|_F^2 \geq \tau$  (i.e.,  $F_i$  is the set of indices of true clusters whose covariances are far from that of the  $i$ -th cluster in Frobenius norm). For every row  $j \in [n]$ , let  $C(j) \in [k]$  be such that  $j \in X_{C(j)}$ . Let's call  $j$ -th row of  $M$  "good" if  $x_j = y_j$  (i.e  $j$ -th sample is not an outlier) and the following condition holds:

$$\sum_{r \in F_{C(j)}} \sum_{\ell \in X_r: x_\ell = y_\ell} M(j, \ell) \leq \|M_j\|_1 \left( \frac{\beta}{\eta} \right).$$

Thus, by Markov's inequality, the fraction of non-outlier entries in  $B_j$  that come from  $X_{r'}$  such that  $r' \in F_r$  is at most  $2 \left( \frac{\beta}{\eta \alpha^2} \right)$ .

Let us estimate the fraction of good rows now. From Lemma 4.6 and Fact 2.18, we have that for every  $r$  and  $r' \in F_r$ :

$$\tilde{\mathbb{E}}[z'(X_r)z'(X_{r'})] \leq \beta.$$

Here, recall that  $z'(X_r) = \frac{1}{w_r n} \sum_{i \leq n} z_i \mathbf{1}(y_i = x_i)$  for every  $r$ . Summing up over  $r' \in F_r$  yields:

$$\frac{1}{w_r n} \sum_{r' \in F_r} \sum_{i \in X_r: x_i = y_i} \sum_{j \in X_{r'}: x_j = y_j} \tilde{\mathbb{E}}[z_i z_j] \leq n \beta.$$

Thus, by Markov's inequality, with probability at least  $1 - \eta$  over the choice of  $i \in X_r$  such that  $x_i = y_i$ , it must hold that:

$$\sum_{r' \in F_r} \sum_{j \in X_{r'}: x_j = y_j} \tilde{\mathbb{E}}[z_i z_j] \leq n \left( \frac{\beta}{\eta} \right).$$

Using that  $(1 - \varepsilon/\alpha)$ -fraction of  $i \in X_r$  satisfy  $x_i = y_i$ , for every  $r$ , we conclude that  $1 - \eta - \varepsilon/\alpha$ -fraction of the rows  $X_r$  are good.

Thus, with probability at least  $(1 - \eta - \varepsilon/\alpha)^\ell \geq (1 - O(\ell(\eta + \varepsilon/\alpha)))$ , every candidate cluster picked in Step 1 of our rounding algorithm corresponds to the large entries from a good row of  $M$ .

We next claim that we cover most of the points in the input in the union of the candidate clusters:

$$|\cup_{i \leq \ell} \hat{C}_i| \geq \left( 1 - 2\sqrt{\eta} - \frac{\varepsilon}{\sqrt{\eta}\alpha} \right) n \tag{7.6}$$

with probability at least  $1 - \sqrt{\eta}$ . To see why, let's estimate the chance that an element  $j \in [n]$  does not appear in any of the  $\hat{C}_i$ s. First, we can assume that  $j$ -th row of  $M$  is good (this loses us  $\eta + \varepsilon/\alpha$ -fraction  $j$ s). For each such  $j$ , there are at least  $\alpha n/2$  large entries. Since  $M$  is symmetric, the  $j$ -th column of  $M$  also has  $\alpha n/2$  large entries. Further,  $j$  appears in  $\cup_{i \leq \ell} \hat{C}_i$  if at least one of the  $\alpha n/2$  large entries are chosen in our rounding. The chance that this does not happen in any of the  $\ell$  picks is at most  $(1 - \alpha/2)^\ell$ . Since  $\ell = \Theta\left(\frac{1}{\alpha} \log(1/\eta)\right)$ , this chance is at most  $O(\eta)$ . Thus, in

expectation  $|\llbracket n \rrbracket \setminus \cup_{i \leq \ell} \hat{C}_i| \leq O(\eta + \varepsilon/\alpha)n$ . By Markov's inequality, with probability at least  $1 - \sqrt{\eta}$ ,  $|\llbracket n \rrbracket \setminus \cup_{i \leq \ell} \hat{C}_i| \leq O\left(\sqrt{\eta} + \frac{\varepsilon}{\sqrt{\eta}\alpha}\right)n$ .

By a union bound, with probability at least  $1 - O(\eta\ell - \varepsilon\ell/\alpha) - \sqrt{\eta} \geq 1 - O(\eta \log(1/\eta)/\alpha - \sqrt{\eta})$ , we must thus have both the following events hold simultaneously:

$$|\cup_{i \leq \ell} \hat{C}_i| \geq \left(1 - 2\sqrt{\eta} - \frac{\varepsilon}{\sqrt{\eta}\alpha}\right)n \geq (1 - 3\sqrt{\eta})n \quad (7.7)$$

and, for every  $1 \leq i \leq \ell$ ,

$$|\hat{C}_i \cap (\cup_{r' \in F_{C(r)}} X_{r'})| \leq 2\left(\frac{\beta}{\eta\alpha} + \varepsilon/\alpha\right) \cdot |\hat{C}_i|. \quad (7.8)$$

**Merging Candidate Clusters** Observe, following the proof of Theorem 4.3, we know that there exists a partition of  $Y$  into sets  $Y_1$  and  $Y_2$  such that for all  $i$ ,

$$\max\left\{\frac{1}{w_i n} |Y_1 \cap X_i|, \frac{1}{w_i n} |Y_2 \cap X_i|\right\} \geq 1 - O(\sqrt{\eta}) - O\left(\frac{\beta}{\eta\alpha^2}\right),$$

and

$$\max_i \frac{1}{w_i n} |X_i \cap Y_1|, \max_i \frac{1}{w_i n} |X_i \cap Y_2| \geq 1 - O(\sqrt{\eta}) - O\left(\frac{\beta}{\eta\alpha^2}\right).$$

Next, we show that the merging step preserves this partition. For each  $\hat{C}_i$ , let  $\hat{C}'_i = \hat{C}_i \cap \cup_{j \notin F_{C(i)}} X_j$ . That is,  $\hat{C}'_i$  is the subset of  $\hat{C}_i$  obtained by removing points from "far-off" clusters and the outliers. Then, since we know that  $|\hat{C}_i| \geq \alpha n/2$  and  $|X \cap Y| \geq (1 - \varepsilon)n$ , we must have  $|\hat{C}_i| - |\hat{C}'_i| = \eta'_i |\hat{C}_i| \leq \left(\frac{3\varepsilon}{\alpha} + \frac{2\beta}{\eta\alpha^2}\right) |\hat{C}_i|$ , where we note that  $\eta'_i \leq \left(\frac{3\varepsilon}{\alpha} + \frac{2\beta}{\eta\alpha^2}\right)$ .

Thus,  $\hat{C}'_i$  is a collection of  $\geq (1 - \eta'_i)\alpha n/2$  points from the submixture  $\cup_{j \notin F_{C(i)}} X_j$ . We know that each  $\mu_i$  is of  $\ell_2$  norm at most  $1/\alpha$ , each  $\Sigma_i$  has spectral norm at most  $1/\alpha$  and that for every  $r, r' \notin F_{C(i)}$ ,  $\|\Sigma_r - \Sigma_{r'}\|_F^2 \leq \tau$ . Further,  $\Sigma_r$  is at most  $\tau + 1/\alpha = O(\tau)$ -different in Frobenius norm from the covariance of the sub-mixture. By an argument similar to the proof of Lemma 2.44, we can establish that the submixture with components  $r$  such that  $r \notin F_{C(i)}$  is  $O(\tau)$ -certifiably bounded variance. Since  $\hat{C}'_i$  is a subset of this sub-mixture of size  $\alpha n/2$ , we immediately obtain that  $\hat{C}'_i$  is  $O(\tau/\alpha)$ -certifiably bounded variance. Thus, applying Lemma 7.7 with outlier parameter  $\eta'_i$  to input  $\hat{C}_i$  yields an estimate  $S_i$  of the 2nd moment of  $\hat{C}'_i$  within a Frobenius error of at most  $O(\tau/\alpha)$ . From Lemma 7.8, this is an additional  $O(1/\alpha)$  different in Frobenius norm from the 2nd moment of the sub-mixture which, as argued above, is itself at most  $O(\tau)$  different in Frobenius norm from  $\Sigma_i$ . Chaining together yields that  $\|\Sigma_i - S_i\|_F^2 \leq O(\tau/\alpha)$  for some constant  $C$ .

Since for every  $r \in S, r' \in T$  it holds that  $\|\Sigma_r - \Sigma_{r'}\|_F^2 \gg \Omega(\tau/\alpha)$ , conditioned on the good event above, our algorithm never merges  $\hat{C}_i$  and  $\hat{C}_j$  whenever  $i, j$  are non-outliers and  $i$  is in some cluster in  $S$  and  $j$  is in some cluster in  $T$ . On the other hand, if  $i, j$  belong to the same cluster, then, the corresponding estimate  $\|S_i - S_j\|_F^2 \leq 2C\tau$ . Thus, our merging process always merges together any such candidates.

As a result, the output of the merging process can have at most one  $i$  from any true cluster – thus, the number of distinct members of  $\mathcal{L}'$  is at most  $k$ . We note that the running time is dominated by computing a pseudo distribution satisfying the union of all the constraints (Step 3 in Algorithm 7.2) and requires  $n^{O(t)}$  time. Step 4 computes a degree  $O(1)$  sos relaxation for at most  $O(\ell)$  components and the merging only requires a fixed polynomial in  $d$  and  $k$  time.  $\square$

## 7.2 Proof of Lemma 7.4

In the following lemma, we show that the constraint system  $\mathcal{A}$ , via a low-degree sum-of-squares proof, implies that a lower bound on the variance of any degree 2 polynomial on  $X'$  whenever the cluster  $\hat{C}$  (indicated by  $z$ ) appreciably intersects two well-separated true clusters.

**Lemma 7.5** (Lower-Bound on Variance of Degree 2 Polynomials). *Let  $Q \in \mathbb{R}^{d \times d}$  be any fixed matrix. Then, for any  $i, j \leq k$ , and  $z'(X_r) = \frac{1}{w_r n} \sum_{i \in X_r} z_i \cdot \mathbf{1}(y_i = x_i)$ , we have for any  $r \neq r' \in [k]$ ,*

$$\mathcal{A} \Big|_{4t}^z \left\{ z'(X_r) z'(X_{r'}) \leq \frac{(32Ct/\alpha)^{2t}}{(\mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q)^{2t}} \left( \frac{\alpha^4}{w_r^2 w_{r'}^2} \left( \mathbb{E}_z (Q - \mathbb{E}_z Q)^2 \right)^t \right. \right. \\ \left. \left. + \frac{\alpha^2}{w_r^2} \left( \mathbb{E}_{X_{r'}} (Q - \mathbb{E}_{X_{r'}} Q)^2 \right)^t + \frac{\alpha^2}{w_{r'}^2} \left( \mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^2 \right)^t \right) \right\}.$$

*Proof.* Let  $z'_i = z_i \mathbf{1}(y_i = x_i)$  for every  $i$ . For every  $1 \leq r \leq k$ , let  $\mathbb{E}_{X_r} Q$  denote the expectation of the homogenous degree 2 polynomial defined by  $Q$ :  $\mathbb{E}_{X_r} Q = \frac{1}{w_r n} \sum_{i, j \in X_r} Q(x_i)$  for every  $r$  where  $Q(x_i) = x_i^\top Q x_i$ . Similarly, let  $\mathbb{E}_z Q$  be the quadratic polynomial in  $z$  defined by  $\mathbb{E}_z Q = \frac{1}{\alpha n} \sum_{i \leq n} z_i Q(x_i)$ . Using the substitution rule and non-negativity of the  $z'_i$ 's, we have for any  $r, r' \in [k]$ :

$$\mathcal{A} \Big|_{4t}^z \left\{ \mathbb{E}_z (Q - \mathbb{E}_z Q)^{2t} = \frac{1}{\alpha n} \sum_{i \in [n]} z_i \left( Q(x_i) - \mathbb{E}_z Q \right)^{2t} \right. \\ \left. \geq \frac{1}{\alpha n} \sum_{i \in X_r \cup X_{r'} : x_i = y_j} z'_i \left( Q(x_i) - \mathbb{E}_z Q \right)^{2t} \right\} \quad (7.9)$$

Then, using the SoS almost triangle inequality (Fact 2.21), we have:

$$\begin{aligned}
\mathcal{A} \Big|_{4t}^z & \left\{ \frac{1}{\alpha n} \sum_{i \in X_r \cup X_{r'}} z'_i \left( Q(x_i) - \mathbb{E}_z Q \right)^{2t} \right. \\
& \geq 2^{-2t} \left( \frac{1}{\alpha n} \sum_{i \in X_r} z'_i \left( \mathbb{E}_{X_r} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha n} \sum_{i \in X_r} z'_i \left( Q(x_i) - \mathbb{E}_{X_r} Q \right)^{2t} \right) \\
& \quad + 2^{-2t} \left( \frac{1}{\alpha n} \sum_{i \in X_r; i: x_i = y_i} z'_i \left( \mathbb{E}_{X_{r'}} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha n} \sum_{i \in X_r; x_i = y_i} z'_i \left( Q(x_i) - \mathbb{E}_{X_{r'}} Q \right)^{2t} \right) \\
& = 2^{-2t} \left( \frac{w_r}{\alpha} z'(X_r) \left( \mathbb{E}_{X_r} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha n} \sum_{i \in X_r} \left( Q(x_i) - \mathbb{E}_{X_r} Q \right)^{2t} \right) \\
& \quad \left. + 2^{-2t} \left( \frac{w_{r'}}{\alpha} z'(X_{r'}) \left( \mathbb{E}_{X_{r'}} Q - \mathbb{E}_z Q \right)^{2t} - \frac{1}{\alpha n} \sum_{i \in X_{r'}} \left( Q(x_i) - \mathbb{E}_{X_{r'}} Q \right)^{2t} \right) \right\} \quad (7.10)
\end{aligned}$$

Next, observe that by the SoS almost triangle inequality (Fact 2.21), we must have:

$$\mathcal{A} \Big|_{4t} \left\{ \left( \mathbb{E}_{X_r} Q - \mathbb{E}_z Q \right)^{2t} + \left( \mathbb{E}_{X_{r'}} Q - \mathbb{E}_z Q \right)^{2t} \geq 2^{-2t} \left( \mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q \right)^{2t} \right\}.$$

Further, note that  $\mathcal{A} \Big|_{O(1)} \left\{ \frac{w_r}{\alpha} z'(X_r) + \frac{w_{r'}}{\alpha} z'(X_{r'}) \leq \frac{1}{\alpha n} \sum_i z_i \leq 1 \right\}$ . Thus, using Fact 4.8 with  $A = \frac{w_r}{\alpha} z'(X_r)$ ,  $B = \frac{w_{r'}}{\alpha} z'(X_{r'})$ ,  $C = (\mathbb{E}_{X_r} Q - \mathbb{E}_z Q)^{2t}$ , and  $D = (\mathbb{E}_{X_{r'}} Q - \mathbb{E}_z Q)^{2t}$  and  $\tau = 2^{-2t} (\mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q)^{2t}$ , we can derive:

$$\begin{aligned}
\mathcal{A} \Big|_{4t}^z & \left\{ (Ct/\alpha)^{2t} \left( \mathbb{E}_z (Q - \mathbb{E}_z Q)^2 \right)^t \geq \mathbb{E}_z (Q - \mathbb{E}_z Q)^{2t} \geq 2^{-6t} \frac{w_r w_{r'}}{\alpha^2} z'(X_r) z'(X_{r'}) \left( \mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q \right)^{2t} \right. \\
& \quad - 2^{-6t} \frac{w_r}{\alpha} \mathbb{E}_{X_r} \left( Q - \mathbb{E}_{X_r} Q \right)^{2t} - 2^{-6t} \frac{w_{r'}}{\alpha} \mathbb{E}_{X_{r'}} \left( Q - \mathbb{E}_{X_{r'}} Q \right)^{2t} \\
& \geq 2^{-6t} \frac{w_r w_{r'}}{\alpha^2} z'(X_r) z'(X_{r'}) (\mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q)^{2t} - \frac{w_r}{\alpha} (Ct/\alpha)^{2t} \left( \mathbb{E}_{X_r} \left( Q - \mathbb{E}_{X_r} Q \right)^2 \right)^t \\
& \quad \left. - \frac{w_{r'}}{\alpha} (Ct/\alpha)^{2t} \left( \mathbb{E}_{X_{r'}} \left( Q - \mathbb{E}_{X_{r'}} Q \right)^2 \right)^t \right\} \quad (7.11)
\end{aligned}$$

where the first inequality uses the Certifiable Hypercontractivity constraints ( $\mathcal{A}_4$ ) and the last inequality follows from the Certifiable Hypercontractivity of  $X_r$  and  $X_{r'}$  (Condition 2.47). Rearranging completes the proof.  $\square$

We can use the lemma above to obtain a simultaneous intersection bound guarantee when there are relative Frobenius separated components in the mixture.

**Lemma 7.6** (Lemma 4.6, restated). *Suppose  $\|\Sigma^{-1/2}(\Sigma_r - \Sigma_{r'})\Sigma^{-1/2}\|_F^2 \geq 10^8 \frac{C^6 t^4}{\beta^{2/t} \alpha^4}$ . Then, for  $z(X_r) = \frac{1}{w_r n} \sum_{i \in X_r} z_i \cdot \mathbf{1}(y_i = x_i)$ ,*

$$\mathcal{A} \Big|_{2t} \{z(X_r)z(X_{r'}) \leq \beta\} .$$

*Proof.* WLOG, we will work with the transformed points  $x_i \rightarrow \Sigma^{-1/2}x_i$  where  $\Sigma$  is the covariance of the mixture. Note that our algorithm does not need to know  $\Sigma$  – this transformation is only for simplifying notation in the analysis that follows.

Let  $\tilde{\Sigma}_z = \Sigma^{-1/2}\Sigma_z\Sigma^{-1/2}$ ,  $\tilde{\Sigma}_r = \Sigma^{-1/2}\Sigma_r\Sigma^{-1/2}$  and  $\tilde{\Sigma}_{r'} = \Sigma^{-1/2}\Sigma_{r'}\Sigma^{-1/2}$  be the transformed covariances. Then, notice that  $\|\tilde{\Sigma}_r\|_2 \leq \frac{1}{w_r} \|\Sigma\|_2 \leq \frac{1.5}{w_r}$  and  $\|\tilde{\Sigma}_{r'}\|_2 \leq \frac{1}{w_{r'}} \|\Sigma\|_2 \leq \frac{1.5}{w_{r'}}$ .

We now apply Lemma 4.10 with  $Q = \tilde{\Sigma}_r - \tilde{\Sigma}_{r'}$ . Then, notice that  $\mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q = \|\tilde{\Sigma}_r - \tilde{\Sigma}_{r'}\|_F^2 + \mu_r^\top (\tilde{\Sigma}_r - \tilde{\Sigma}_{r'}) \mu_r - \mu_{r'}^\top (\tilde{\Sigma}_r - \tilde{\Sigma}_{r'}) \mu_{r'} \geq \|Q\|_F^2 - \frac{4}{\alpha}$ . Then, we obtain:

$$\begin{aligned} & \mathcal{A} \Big|_{2t} \left\{ z(X_r)z(X_{r'}) \right. \\ & \leq \left( \frac{32Ct/\alpha}{\mathbb{E}_{X_r} Q - \mathbb{E}_{X_{r'}} Q} \right)^{2t} \left( \frac{\alpha^2}{w_r w_{r'}} \left( \mathbb{E}_z (Q - \mathbb{E}_z Q)^2 \right)^t + \frac{\alpha}{w_r} \left( \mathbb{E}_{X_{r'}} (Q - \mathbb{E}_{X_{r'}} Q)^2 \right)^t + \frac{\alpha}{w_{r'}} \left( \mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^2 \right)^t \right) \left. \right\} . \end{aligned} \tag{7.12}$$

Since  $X_r$  and  $X_{r'}$  have certifiably  $C$ -bounded variance polynomials for  $C = 4$  (as a consequence of Condition 2.47 and Fact 2.45 followed by an application of Lemma 2.25), we have:

$$\mathbb{E}_{X_{r'}} (Q - \mathbb{E}_{X_{r'}} Q)^2 \leq 6 \left\| \tilde{\Sigma}_{r'}^{1/2} Q \tilde{\Sigma}_{r'}^{1/2} \right\|_F^2 \leq \frac{10}{w_{r'}^2} \|Q\|_F^2 \leq \frac{10}{\alpha^2} \|Q\|_F^2 ,$$

and

$$\mathbb{E}_{X_r} (Q - \mathbb{E}_{X_r} Q)^2 \leq 6 \left\| \tilde{\Sigma}_r^{1/2} Q \tilde{\Sigma}_r^{1/2} \right\|_F^2 \leq \frac{10}{w_r^2} \|Q\|_F^2 \leq \frac{10}{\alpha^2} \|Q\|_F^2 .$$

Finally, using the bounded-variance constraints in  $\mathcal{A}$ , we have:

$$\mathcal{A} \Big|_{\frac{Q,z}{4}} \mathbb{E}(Q - \mathbb{E}_z Q)^2 \leq \frac{10}{\alpha^2} \|Q\|_F^2 .$$

Plugging these estimates back in (7.12) yields:

$$\mathcal{A} \Big|_{\frac{z}{4}} \left\{ z(X_r)z(X_{r'}) \leq \frac{(1000Ct/\alpha)^{2t}}{\|Q\|_F^{2t} \alpha^{2t}} \left( \alpha^2 + \frac{\alpha}{w_r} + \frac{\alpha}{w_r w_{r'}} \right) \leq \frac{3}{w_r w_{r'}} \frac{(1000Ct)^{2t}}{\alpha^{2t} \|Q\|_F^{2t}} \leq \frac{3(1000Ct)^{2t}}{\alpha^{2t} \|Q\|_F^{2t}} \right\} . \tag{7.13}$$

Plugging in the lower bound on  $\|Q\|_F^{2t}$  and applying cancellation within SoS (Fact 4.9) completes the proof.  $\square$

### 7.3 2nd Moment Estimation Subroutine

The following lemma gives a 2nd moment estimation algorithm with error in Frobenius norm for distributions that have a certifiably bounded covariance. The proof is very similar to the SoS based mean and covariance estimation algorithms but we provide it in full for completeness here.

**Lemma 7.7** (2nd Moment Estimation in Frobenius Norm). *Let  $1/100 \geq \eta > 0$ . There is an  $n^{O(1)}$  time algorithm that takes input an  $\eta$ -corruption  $Y$  of an sample  $X$  of size  $n$  and outputs an estimate  $M_2$  of the 2nd moment of  $X$  with the following properties: Let  $X \subseteq \mathbb{R}^d$  be a collection of  $n$  points satisfying  $\left| \frac{Q}{2} \left\{ \frac{1}{|X|} \sum_{x \in X} \left( Q(x) - \frac{1}{|X|} Q(x) \right)^2 \leq C \|Q\|_F^2 \right\} \right.$  for a matrix-valued indeterminate  $Q$ . Let  $M_2 = \frac{1}{n} \sum_{x \in X} x x^\top$ . Then, the estimate  $\hat{M}_2$  output by the algorithm satisfies:*

$$\left\| \hat{M}_2 - M_2 \right\|_F^2 \leq 80C\eta.$$

*Proof.* Consider the constraint system with scalar-valued indeterminates  $z_i$  for  $1 \leq i \leq n$  and  $d$ -dimensional vector-valued indeterminates  $x'_1, x'_2, \dots, x'_n$  with the following set of constraints:  $\mathcal{A} =$

$$\left. \begin{array}{l} \forall i \leq n \\ \\ \\ \forall i \leq n \\ \frac{1}{n} \sum_{i=1}^n \left( x_i'^\top Q x_i - \frac{1}{n} \sum_{i=1}^n x_i'^\top Q x_i' \right)^2 \leq C \|Q\|_F^2 \end{array} \right\} \begin{array}{l} z_i^2 = z_i \\ \sum_{i=1}^n z_i = (1 - \eta)n \\ \tilde{M}_2 = \frac{1}{n} \sum_{i=1}^n x_i' x_i'^\top \\ z_i x_i' = z_i y_i \end{array} \quad (7.14)$$

Observe that  $X' = X$  and  $z_i$  set to the 0-1 indicator of non-outliers satisfies the constraint system. Thus, the constraints are feasible.

Our algorithm finds a pseudo-distribution  $\tilde{\zeta}$  of degree 10 satisfying the above constraints and output  $\tilde{\mathbb{E}}[\tilde{M}_2]$ . Let us now analyze this algorithm. The key is the following statement that gives a sum-of-squares proof of closeness of  $\tilde{M}_2$  and  $M_2$  in Frobenius norm. We use the notation  $\mathbb{E}_X Q$  and  $\mathbb{E}_{X'} Q$  to abbreviate  $\frac{1}{n} \sum_{i=1}^n x_i^\top Q x_i$  and  $\frac{1}{n} \sum_{i=1}^n x_i'^\top Q x_i'$  respectively.

$$\begin{aligned} \mathcal{A} \Big|_{\frac{Q}{2}} & \left\{ \left( \frac{1}{n} \sum_{i=1}^n x_i^\top Q x_i - \frac{1}{n} \sum_{i=1}^n x_i'^\top Q x_i' \right)^2 \right. \\ & = \left( \frac{1}{n} \sum_{i=1}^n (1 - z_i \mathbf{1}(x_i = y_i)) x_i^\top Q x_i - x_i'^\top Q x_i' \right)^2 \\ & \leq \left( \frac{1}{n} (1 - z_i \mathbf{1}(x_i = y_i))^2 \right) \left( \frac{1}{n} \sum_{i=1}^n (x_i^\top Q x_i - x_i'^\top Q x_i')^2 \right) \end{aligned}$$

$$\begin{aligned}
&\leq 20\eta \cdot \left( \frac{1}{n} \sum_{i=1}^n \left( x_i^\top Q x_i - \mathbb{E}_X Q \right)^2 + \frac{1}{n} \sum_{i=1}^n \left( x'_i{}^\top Q x'_i - \mathbb{E}_{X'} Q \right)^2 + \left( \mathbb{E}_X Q - \mathbb{E}_{X'} Q \right)^2 \right) \\
&\leq 20\eta(2C \|Q\|_F^2) + 20\eta \left( \mathbb{E}_X Q - \mathbb{E}_{X'} Q \right)^2 \Big\}
\end{aligned}$$

where the first inequality follows by the SoS version of the Cauchy-Schwarz inequality and the 2nd by the SoS version of the Almost Triangle inequality.

Rearranging and using that  $1 - 20\eta > 1/2$  now yields that:

$$\mathcal{A} \Big|_{\frac{Q}{2}} \left\{ \left( \frac{1}{n} \sum_{i=1}^n x_i^\top Q x_i - \frac{1}{n} \sum_{i=1}^n x'_i{}^\top Q x'_i \right)^2 \leq 80C\eta \|Q\|_F^2 \right\}$$

Notice that the LHS above equals the linear polynomial  $\langle \tilde{M}_2 - M_2, Q \rangle$ . We now plug in  $Q = \tilde{M}_2 - M_2$  to obtain:

$$\mathcal{A} \Big|_{\frac{Q}{2}} \left\{ \|\tilde{M}_2 - M_2\|_F^4 \leq 80C\eta \|\tilde{M}_2 - M_2\|_F^2 \right\}$$

Applying Fact 2.24 yields:

$$\mathcal{A} \Big|_{\frac{Q}{2}} \left\{ \|\tilde{M}_2 - M_2\|_F^8 \leq 80^4 C^4 \eta^4 \right\}$$

Taking pseudo-expectations with respect to  $\tilde{\zeta}$  and using Hölder's inequality for pseudo-distributions yields that

$$\left\| \tilde{\mathbb{E}}_{\tilde{\zeta}} \tilde{M}_2 - M_2 \right\|_F^8 \leq \tilde{\mathbb{E}}_{\tilde{\zeta}} \|\tilde{M}_2 - M_2\|_F^8 \leq 80^4 C^4 \eta^4.$$

Taking the 4-th root, we can conclude our rounded value  $\hat{M}_2 = \tilde{\mathbb{E}}_{\tilde{\zeta}} \tilde{M}_2$  satisfies:

$$\left\| \hat{M}_2 - M_2 \right\|_F^2 \leq 80C\eta.$$

This completes the proof. □

We also note the following simple consequence of the certifiable bounded variance property that follows via an argument similar to the one employed in the proof of the previous lemma.

**Lemma 7.8** (Subsamples of Bounded-Variance Distributions). *Let  $X \subseteq \mathbb{R}^d$  be a collection of  $n$  points satisfying  $\left| \frac{Q}{2} \left\{ \frac{1}{|X|} \sum_{x \in X} (Q(x) - \frac{1}{|X|} Q(x))^2 \leq C \|Q\|_F^2 \right\} \right.$  for a matrix-valued indeterminate  $Q$ . Let  $M_2 = \frac{1}{|X|} \sum_{x \in X} x x^\top$  be the 2nd moment of  $X$ . Let  $S \subseteq X$  be a subset of size at least  $\beta|X|$ . Then,*

$$\left| \frac{1}{4} \left\{ \left\| \frac{1}{|S|} \sum_{x \in S} x x^\top - M_2 \right\|_F^2 \leq \frac{1}{\beta} \right\} \right.$$



*Proof.* We have by the Cauchy-Schwarz inequality:

$$\begin{aligned} \frac{|Q|}{2} \left\{ \left( \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(x \in S)(x^\top Qx - M_2), Q \right) \right\}^2 &\leq \left( \frac{1}{|X|} \mathbf{1}(x \in S)^2 \right) \left( \sum_{x \in X} (x^\top Qx - M_2), Q \right)^2 \\ &\leq \left( \frac{|S|}{|X|} \right) \|Q\|_F^2. \end{aligned}$$

We now substitute in  $Q = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(x \in S)(x^\top Qx - M_2)$  to obtain:

$$\frac{|Q|}{2} \left\| \left\| \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(x \in S)(x^\top Qx - M_2) \right\|_F \right\|^4 \leq \left( \frac{|S|}{|X|} \right) \left\| \left\| \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(x \in S)(x^\top Qx - M_2) \right\|_F \right\|^2.$$

We now apply Fact 2.24 to yield:

$$\frac{|Q|}{2} \left\| \left\| \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(x \in S)(x^\top Qx - M_2) \right\|_F \right\|^8 \leq \left( \frac{|S|}{|X|} \right)^4.$$

We finally apply Fact 4.9 to conclude that:

$$\frac{|Q|}{2} \left\| \left\| \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(x \in S)(x^\top Qx - M_2) \right\|_F \right\|^2 \leq \left( \frac{|S|}{|X|} \right).$$

Rescaling gives the claim. □

## 8 Getting $\text{poly}(\varepsilon)$ -close in TV Distance: Proof of Theorem 1.5

**Theorem 8.1** (Robustly Learning  $k$ -Mixtures with small error). *Given  $0 < \varepsilon < 1/k^{k^{O(k^2)}}$  and a multiset  $Y = \{y_1, y_2, \dots, y_n\}$  of  $n$  i.i.d. samples from a distribution  $F$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , for an unknown  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_{i \leq k} w_i \mathcal{N}(\mu_i, \Sigma_i)$ , where  $n \geq n_0 = d^{O(k)} \text{poly}_k(1/\varepsilon)$ , there exists an algorithm that runs in time  $n^{O(1)} \text{poly}_k(1/\varepsilon)$  and with probability at least 0.99 outputs a hypothesis  $k$ -mixture of Gaussians  $\widehat{\mathcal{M}} = \sum_{i \leq k} \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  such that  $d_{\text{TV}}(\mathcal{M}, \widehat{\mathcal{M}}) = \mathcal{O}(\varepsilon^{c_k})$ , with  $c_k = 1/(100^k C^{(k+1)!} k! \text{sf}(k+1))$ , where  $C > 0$  is a universal constant and  $\text{sf}(k) = \prod_{i \in [k]} (k-i)!$  is the super-factorial function.*

In order to obtain the above theorem, we require recovering a polynomial sized list of candidate parameters, in addition to the efficient partial clustering result we obtained in the previous section. To this end, we show the following list-recovery theorem which is similar to Theorem 3.1, but the algorithm outputs a polynomial-size list instead.

**Theorem 8.2** (Recovering a small list of candidate parameters). *Fix any  $\alpha > \varepsilon > 0, \Delta > 0$ . Let  $X$ , a sample from a  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  satisfying Condition 2.47 with parameters  $\gamma = \varepsilon d^{-8k} k^{-Ck}$ , for  $C$  a sufficiently large universal constant, and  $t = 8k$ , and let  $Y$  be an  $\varepsilon$ -corruption of  $X$ . Let  $X'$  be a set of  $n' = \mathcal{O}(\varepsilon \eta / (k^5 (\Delta^4 + 1/\alpha^4)))^{-4ck}$  fresh samples from  $\mathcal{M}$  and  $Z$  be an  $\varepsilon$ -corruption of*

$X'$ . If  $w_i \geq \alpha$ ,  $\|\mu_i\|_2 \leq \frac{2}{\sqrt{\alpha}}$  and  $\|\Sigma_i - I\|_F \leq \Delta$  for every  $i \in [k]$ , then, given  $k, Y, Z$  and  $\varepsilon$ , the algorithm outputs a list  $L$  of at most  $\ell' = O\left(\left(k^5 (\Delta^4 + 1/\alpha^4)\right)^{4k} / \eta^{4k}\right)$  candidate hypotheses (component means and covariances), such that with probability at least  $99/100$  there exist  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]} \subseteq L$  satisfying  $\|\mu_i - \hat{\mu}_i\|_2 \leq O\left(\frac{\Delta^{1/2}}{\alpha}\right) \eta^{G(k)}$  and  $\left\|\Sigma_i - \hat{\Sigma}_i\right\|_F \leq O(k^4) \frac{\Delta^{1/2}}{\alpha} \eta^{G(k)}$ , for all  $i \in [k]$ . Here,  $\eta = (2k)^{4k} O(1/\alpha + \Delta)^{4k} \varepsilon^{1/(kO(k^2))}$  and  $G(k) = \frac{1}{C^{k+1}(k+1)!}$ . The running time of the algorithm is  $\text{poly}(|L|, |Y|, d^k) \cdot \text{poly}_k(1/\varepsilon)$ .

## 8.1 Proof of Theorem 8.2

We use the following notation and background from Moitra-Valiant [MV10]:

**Definition 8.3** (Statistically Learnable). Given  $\varepsilon > 0$ , we call a mixture of Gaussians  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$   $\varepsilon$ -statistically learnable if  $\min_i w_i \geq \varepsilon$  and  $\min_{i \neq j} d_{TV}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j)) \geq \varepsilon$ .

**Definition 8.4** (Correct Subdivision). Given a Gaussian mixture of  $k$  Gaussians,  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  and a mixture of  $k' \leq k$  Gaussians  $\hat{\mathcal{M}} = \sum_i \hat{w}_i \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i)$ , we call  $\hat{\mathcal{M}}$  an  $\varepsilon$ -correct subdivision of  $\mathcal{M}$  if there is a function  $\pi : [k] \rightarrow [k']$  that is onto and

1.  $\forall j \in [k'], \left| \sum_{i: \pi(i)=j} w_i - \hat{w}_j \right| \leq \varepsilon$
2.  $\forall i \in [k], \left\| \mu_i - \hat{\mu}_{\pi(i)} \right\| + \left\| \Sigma_i - \hat{\Sigma}_{\pi(i)} \right\|_F \leq \varepsilon$ .

**Theorem 8.5** (Theorem 8 in [MV10]). Given an  $\varepsilon$ -statistically learnable Gaussian mixture  $\mathcal{M}$  in isotropic position, for some  $\varepsilon > 0$ , there exists an algorithm that requires  $n = \text{poly}(d/\varepsilon)$  samples and runs in time  $O(\text{poly}_k(n))$  and with probability at least  $99/100$  recovers an  $\varepsilon$ -correct sub-division  $\hat{\mathcal{M}}$ . Let the corresponding algorithm be referred to as *PARTITION PURSUIT*.

The algorithm has two steps: first run the first three steps of Algorithm 3.2 to get the list  $L'$  of  $\hat{S}$  and  $V'_S$ ; then apply the following proposition to learn the mixture in the subspace  $V'_S$ . This proposition is a generalization of Theorem 8.5 without the assumption that the total variation distance between each pair of components is at least  $\varepsilon$ . The sample and time complexities has a worse, but still polynomial dependence on  $\varepsilon$ . Note that although the algorithm in the proposition is non-robust, we can take a sample without noise with constant probability because the algorithm only requires a polynomial number of samples in  $\varepsilon$ .

**Algorithm 8.6** (Efficient List-Recovery of Candidate Parameters).

**Input:** An  $\varepsilon$ -corruption  $Y$  of a sample  $X$  from a  $k$ -mixture of Gaussians  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$ . Let  $Z$  be an additional  $\varepsilon$ -corrupted sample of size  $n'$  from  $\mathcal{M}$ .

**Requirements:** The guarantees of the algorithm hold if the mixture parameters and the sample  $X$  satisfy:

1.  $w_i \geq \alpha$  for all  $i \in [k]$ ,
2.  $\|\mu_i\|_2 \leq 2/\sqrt{\alpha}$  for all  $i \in [k]$ ,

3.  $\|\Sigma_i - I\|_F \leq \Delta$  for all  $i \in [k]$ .
4.  $X$  satisfies Condition 2.47 with parameters  $(\gamma, t)$ , where  $\gamma = \varepsilon d^{-8k} k^{-Ck}$ , for  $C$  a sufficiently large universal constant, and  $t = 8k$ .
5. The number of fresh samples  $n' = O(\varepsilon \eta / (k^5 (\Delta^4 + 1/\alpha^4)))^{-4ck}$ , for a fixed constant  $c$ .

**Parameters:**  $\eta = (2k)^{4k} O(1/\alpha + \Delta)^{4k} \varepsilon^{1/(k^{O(k^2)})}$ ,  $D = C(k^4/(\alpha\sqrt{\eta}))$ ,  $\delta = 2\eta^{1/(C^{k+1}(k+1)!)}$ ,  $\ell' = 100 \log k (\eta / (k^5 (\Delta^4 + 1/\alpha^4)))^{-4k}$ , for some sufficiently large absolute constant  $C > 0$ ,  $\lambda = 4\eta$ ,  $\phi = 10(1 + \Delta^2)/(\sqrt{\eta}\alpha^5)$ ,  $\varepsilon_1 = O(\sqrt{\Delta}\delta^{1/4}/\alpha)$ .

**Output:** A list  $L$  of hypotheses such that there exists at least one,  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \leq k} \in L$ , satisfying:  
 $\|\mu_i - \hat{\mu}_i\|_2 \leq O\left(\frac{\Delta^{1/2}}{\alpha}\right) \eta^{G(k)}$  and  $\|\Sigma_i - \hat{\Sigma}_i\|_F \leq O(k^4) \frac{\Delta^{1/2}}{\alpha} \eta^{G(k)}$ , where  $G(k) = \frac{1}{C^{k+1}(k+1)!}$ .

**Operation:**

1. **Robust Estimation of Hermite Tensors:** For  $m \in [4k]$ , compute  $\hat{T}_m$  such that  $\max_{m \in [4k]} \|\hat{T}_m - \mathbb{E}[h_m(\mathcal{M})]\|_F \leq \eta$  using the robust mean estimation algorithm in Fact 2.36.
2. **Random Collapsing of Two Modes of  $\hat{T}_4$ :** Let  $L'$  be an empty list. Repeat  $\ell'$  times: For  $j \in [4k]$ , choose independent standard Gaussians in  $\mathbb{R}^d$ , denoted by  $x^{(j)}, y^{(j)} \sim \mathcal{N}(0, I)$ , and uniform draws  $a_1, a_2, \dots, a_t$  from  $[-D, D]$ . Let  $\hat{S}$  be a  $d \times d$  matrix such that for all  $r, s \in [d]$ ,  $\hat{S}(r, s) = \sum_{j \in [4k]} a_j \hat{T}_4(r, s, x^{(j)}, y^{(j)}) = \sum_{j \in [4k]} a_j \sum_{g, h \in [d]} \hat{T}_4(r, s, g, h) x^{(j)}(g) y^{(j)}(h)$ . Add  $\hat{S}$  to the list  $L'$ .
3. **Construct Low-Dimensional Subspace:** Let  $V$  be the span of all singular vectors of the natural  $d \times d^{m-1}$  flattening of  $\hat{T}_m$  with singular values  $\geq \lambda$  for  $m \leq 4k$ . For each  $\hat{S} \in L'$ , let  $V'_\hat{S}$  be the span of  $V$  plus all the singular vectors of  $\hat{S}$  with singular value larger than  $\delta^{1/4}$ .
4. **Moitra-Valiant for Low-Dimensional Subspace:** Initialize  $L$  to be the empty list. For each  $\hat{S} \in L'$ , let  $\hat{P} = UU^\top$  be the orthogonal projection matrix onto the span of  $V'_\hat{S}$ , where  $U \in \mathbb{R}^{d \times d}$  has orthonormal columns. Let  $m = \dim(V'_\hat{S})$  and let  $\hat{Z} \subset Z$  be a randomly chosen subset of size  $\text{poly}(m/\varepsilon_1)$ . Let  $U_m$  denote the first  $m$  columns of  $U$  and for all  $z \in \hat{Z}$ , compute  $U_m^\top z$ . Run PARTITION PURSUIT on the resulting set of points and let  $\{\hat{\mu}_i^{\hat{P}}, \hat{\Sigma}_i^{\hat{P}}\}_{i \in [k]}$  be the parameters corresponding to the  $\varepsilon$ -correct subdivision output by PARTITION PURSUIT. Let  $\hat{\mu}_i^\top = [(\hat{\mu}_i^{\hat{P}})^\top, 0]$  be a  $d$  dimensional vector padded with 0s and  $\hat{\Sigma}_i$  be a  $d \times d$  matrix with  $\hat{\Sigma}_i^{\hat{P}}$  in the top left  $m \times m$  sub-matrix and 0's elsewhere. Add  $\{U\hat{\mu}_i, U\hat{\Sigma}_i U^\top + (\hat{S} + I) - \hat{P}(I + S)\hat{P}\}_{i \in [k]}$  to  $L$ .

**Proposition 8.7.** Given  $\varepsilon > 0$  and a sample  $X$  of size  $\text{poly}(d, 1/\varepsilon)$  from a  $k$ -mixture of Gaussians  $\mathcal{M}$  with mixture covariance  $\Sigma$  such that  $0.99I \leq \Sigma \leq 1.01I$  and satisfies  $w_i \geq \varepsilon$ , the PARTITION PURSUIT algorithm runs in time  $\text{poly}(d, 1/\varepsilon)$  and with probability at least  $9/10$  returns an  $O(\varepsilon)$ -correct sub-division, denoted by  $\hat{\mathcal{M}}$ .

Recall, the PARTITION PURSUIT algorithm satisfies Theorem 8.5 and we will prove that with

an appropriately chosen parameter  $\varepsilon$ , PARTITION PURSUIT also satisfies Proposition 8.7. The main idea is that if any two components are actually close enough in total variation distance, then any algorithm with access to only a polynomial number of samples could never distinguish these two components from a single Gaussian. So if all pairwise distances are either sufficiently large or sufficiently small, the algorithm will behave as if it were given sample access to a mixture that meets the requirements of Theorem 8.5.

**Lemma 8.8.** *Given  $0 < \gamma, \delta < 1$  and two distributions,  $\mathcal{D}_1$  and  $\mathcal{D}_2$  over  $\mathbb{R}^d$  such that  $d_{TV}(\mathcal{D}_1, \mathcal{D}_2) < \gamma$ , let  $X_1$  be set of  $n$  i.i.d. samples from  $\mathcal{D}_1$  and  $X_2$  be  $n$  i.i.d. samples from  $\mathcal{D}_2$ . Let  $\mathcal{A}$  be any algorithm that takes as input  $X_1$  and outputs a list of  $m$  real numbers,  $Y_1 = \{y_i\}_{i \in [m]}$ , such that  $y_i \in [-1, 1]$  with probability at least  $1 - \delta$ . Then, for any  $\tau > 0$ ,  $\mathcal{A}$  on input  $X_2$  outputs a list of  $m$  real numbers  $Y_2 = \{y'_i\}_{i \in [m]}$  such that with probability at least  $1 - \delta - (4mn\gamma/\tau)$ , for all  $i \in [m]$ ,  $|y_i - y'_i| \leq \tau$ .*

*Proof.* Let  $\mathcal{U}_1$  be the uniform distribution over  $X_1$  and  $\mathcal{U}_2$  be the uniform distribution over  $X_2$ . Then,

$$\begin{aligned} d_{TV}(\mathcal{U}_1, \mathcal{U}_2) &\leq \sqrt{2}H^2(\mathcal{U}_1, \mathcal{U}_2) = \sqrt{2}nH^2(\mathcal{D}_1, \mathcal{D}_2) \\ &\leq \sqrt{2}nd_{TV}(\mathcal{D}_1, \mathcal{D}_2) \\ &\leq \sqrt{2}\gamma n \end{aligned} \tag{8.1}$$

Consider the family of functions  $\mathcal{F}$  that take as input  $n$  samples and output a single bit in  $\{0, 1\}$ . We know that for any function  $f \in \mathcal{F}$ , the probability that  $f(X_1) \neq f(X_2)$  is at most  $\sqrt{2}\gamma n$ . Recall, the algorithm outputs  $m$  real numbers in the range  $[-1, 1]$ , which we can discretize into a grid  $\Delta$  of length  $\tau$ . There are at most  $2/\tau$  distinct grid points and for any  $y_i \in [-1, 1]$ , there exists a point  $z_i \in \Delta$  such that  $|y_i - z_i| \leq \tau$ . Further, observe we can represent each  $y_i$  using  $2/\tau$  functions  $f \in \mathcal{F}$ . Then, union bounding over the events that each of the  $2/\tau$  functions output different bits, for each of the  $m$  parameters, we have that with probability at least  $1 - (2\sqrt{2}\gamma nm/\tau)$ , any algorithm outputs a list  $\{y'_i\}_{i \in [m]}$  such that  $|y_i - y'_i| \leq \tau$ . Finally, union bounding over the event that algorithm  $\mathcal{A}$  fails with probability  $\delta$  yields the claim.  $\square$

We then prove there is a gap  $[f(d, \varepsilon_1), \varepsilon_1)$  between pairwise distances of components so that if we merge components within distance  $f(d, \varepsilon_1)$ , the resulting mixture is  $\varepsilon_1$ -statistically learnable.

**Lemma 8.9.** *Let  $f(d)(\varepsilon) = f(d, \varepsilon)$ . There exists  $\ell \in [k^2]$  such that for every pair of components, either  $d_{TV}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j)) < (f(d))^\ell(\varepsilon)$  or  $d_{TV}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j)) \geq (f(d))^{\ell-1}(\varepsilon)$ . Moreover, the set of Gaussians with total variation distance at most  $(f(d))^\ell(\varepsilon)$  is an equivalence class.*

*Proof.* We can see that intervals  $\{(f(d))^\ell(\varepsilon), (f(d))^{\ell-1}(\varepsilon)\}_{\ell \in [k^2]}$  are disjoint. There are at most  $k^2 - 1$  distinct values of  $d_{TV}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j))$ . So there exists an interval  $[(f(d))^\ell(\varepsilon), (f(d))^{\ell-1}(\varepsilon))$  such that for every pair of components  $\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j)$ , either  $d_{TV}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j)) < (f(d))^\ell(\varepsilon)$  or  $d_{TV}(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j)) \geq (f(d))^{\ell-1}(\varepsilon)$ .

Next, we show for any  $\ell$ , Gaussians with pair wise TV distance  $(f(d))^\ell(\varepsilon)$  form an equivalence class. Consider component Gaussians  $G_1, G_2$  and  $G_3$  such that  $G_1$  and  $G_2$  are at total variation

distance at most  $(f(d))^\ell(\varepsilon)$  and  $G_2$  and  $G_3$  are also at total variation distance at most  $(f(d))^\ell(\varepsilon)$ .

$$\begin{aligned} d_{TV}(G_1, G_3) &\leq d_{TV}(G_1, G_2) + d_{TV}(G_2, G_3) \\ &\leq 2(f(d))^\ell(\varepsilon) \\ &\ll (f(d))^{\ell-1}(\varepsilon) \end{aligned}$$

and since there is no pair of Gaussians with total variation distance inside the interval  $[(f(d))^\ell(\varepsilon), (f(d))^{\ell-1}(\varepsilon)]$ , this implies  $d_{TV}(G_1, G_3) \leq (f(d))^\ell(\varepsilon)$ .  $\square$

We can now complete the proof of Proposition 8.7 :

*Proof of Proposition 8.7.* By Lemma 8.9, there exists an interval  $[(f(d))^\ell(\varepsilon), (f(d))^{\ell-1}(\varepsilon)]$  such that there is no pair of Gaussians with total variation distance inside the interval and  $(f(d))^\ell(\varepsilon), (f(d))^{\ell-1}(\varepsilon)$  are polynomials in  $d$  and  $\varepsilon$ . Let  $\varepsilon_1 = (f(d))^{\ell-1}(\varepsilon)$  and  $f(d, \varepsilon_1) = (f(d))^\ell(\varepsilon)$ . Let  $X$  be a set of  $n = (d/\varepsilon)^c$  samples from  $\mathcal{M}$ , where  $c$  is fixed universal constant. Let  $\bar{\mathcal{M}}$  be the mixture obtained by merging all components in an equivalence class with total variation distance at most  $f(d, \varepsilon_1)$  to a single Gaussian and observe  $d_{TV}(\mathcal{M}, \bar{\mathcal{M}}) \leq kf(d, \varepsilon_1)$ . Next, observe that PARTITION PURSUIT outputs at most  $k$  means and covariances, which can be represented as a list of at most  $2kd^2$  real numbers. Further, since  $\Sigma \leq 1.01I$  and  $w_i \geq \varepsilon$ , the means of each component  $\|\mu_i\|_2^2 \leq 2/\varepsilon$  and  $\|\Sigma_i\|_F^2 \leq O(d^2/\varepsilon)$ .

Then, rescaling the instance by  $O(\varepsilon/d^2)$  and applying Lemma 8.8 with  $\mathcal{D}_1 = \mathcal{M}$ ,  $\mathcal{D}_2 = \bar{\mathcal{M}}$ , input samples  $X$  and accuracy parameter  $\tau = (\varepsilon/d)^{c_2}$ , for a large enough constant  $c_2$ , it follows that with probability at least  $1 - 0.99 - O(f(d, \varepsilon_1) \cdot (\varepsilon/d)^{c_3})$ , for a fixed constant  $c_3$ , the resulting list of numbers is  $\tau$ -close to that obtained by running PARTITION PURSUIT on a set of  $n$  samples from  $\bar{\mathcal{M}}$ . Since  $\bar{\mathcal{M}}$  is  $\varepsilon_1$ -statistically learnable, it follows from Theorem 8.5 that with probability at least  $9/10$ , PARTITION PURSUIT will output an  $O(\varepsilon_1)$ -correct sub-division  $\hat{\mathcal{M}}$ .  $\square$

*Proof of Theorem 8.2.* Recall, by part (1) of Proposition 3.5, the dimension of the subspace  $V'_\zeta$  is  $m = \dim(V'_\zeta) = O\left(\frac{(k(1+\Delta+1/\alpha))^{4k+5}}{\eta^2}\right)$  and let  $\varepsilon_1 = \sqrt{\Delta}\delta^{1/4}/\alpha$ . Let  $c_{mv}$  be a fixed constant such that  $(m/\varepsilon_1)^{c_{mv}}$  samples suffice for applying Theorem 8.5. Further, observe in the fresh sample  $Y$ , the probability that any given sample is corrupted is  $\varepsilon$ . Let  $\zeta$  be the event that a random subset of  $(m/\varepsilon_1)^{c_{mv}}$  samples from  $Z$  does not contain any corrupted points. Then, the event  $\zeta$  holds with probability at least  $(1 - \varepsilon)^{(m/\varepsilon_1)^{c_{mv}}}$ . Conditioning on  $\zeta$  and running step 4 of Algorithm 8.6, it follows from Proposition 8.7 that we recover  $O(\varepsilon_1)$ -accurate estimates to the parameters of  $\mathcal{M}$  in the subspace, i.e.  $\|U^\top \mu_i - \hat{\mu}_i\|_2 \leq O(\varepsilon_1)$  and  $\|U^\top \Sigma_i U - \hat{\Sigma}_i\|_F \leq O(\varepsilon_1)$ . Since we repeat the above for  $\ell'$  candidate subspaces in  $L'$ , the probability over all probability of success is  $(1 - \varepsilon)^{(m/\varepsilon_1)^{c_{mv}} \cdot \ell'}$ .

By part (2) in Proposition 3.5, there is a vector  $\mu'_i \in V'_\zeta$  such that  $\|\mu_i - \mu'_i\| \leq \frac{20}{\alpha} \delta^{1/4} \Delta^{1/2}$  where  $\delta = 2\eta^{1/(C^{k+1}(k+1)!)}$  and  $\eta = O(4k(1 + 1/\alpha + \Delta)^{4k} \sqrt{\varepsilon_1})$ . Let  $\hat{P} = UU^\top$  be a projection matrix where the columns of  $Q$  span  $V'_\zeta$  and let  $Q^\top \mu_i$  be the projection of the true means to the corresponding

subspace. Then,

$$\begin{aligned}
\|U\hat{\mu}_i - \mu_i\|_2 &\leq \|U\hat{\mu}_i - \mu'_i\|_2 + \|\mu'_i - \mu_i\|_2 \\
&\leq \|U\hat{\mu}_i - P(\mu'_i - \mu_i + \mu_i)\|_2 + O\left(\frac{\sqrt{\Delta}\delta^{1/4}}{\alpha}\right) \\
&\leq \|\hat{\mu}_i - U^\top \mu_i\|_2 + O\left(\frac{\sqrt{\Delta}\delta^{1/4}}{\alpha}\right) \\
&\leq O\left(\frac{\sqrt{\Delta}\delta^{1/4}}{\alpha}\right).
\end{aligned}$$

where the third inequality follows from observing that  $U^\top \mu_i$  is the true mean in the low dimensional subspace and applying Proposition 8.7.

By Proposition 3.4, there exists  $\hat{S} \in L'$  such that  $\hat{S} - (\Sigma_i - I) = P_i + Q_i$  where  $\|P_i\|_F = O(\sqrt{\eta/\alpha})$ . Again by part (3) in Proposition 3.5, there exists a symmetric matrix  $Q'_i \in V'_{\hat{S}} \times V'_{\hat{S}}$  such that  $\|Q_i - Q'_i\|_F \leq O(\frac{k^2}{\alpha}\delta^{1/4}\Delta^{1/2})$ . We also know that in the subspace spanned by  $V'_{\hat{S}}$ ,  $\|\hat{\Sigma}_i - U^\top \Sigma_i U\|_F^2 \leq \text{poly}(\varepsilon_2)$ . Recall, Algorithm 8.6 outputs the following estimate:  $\hat{M} = U\hat{\Sigma}U^\top + (I + \hat{S}) - \hat{P}(I + \hat{S})\hat{P}$ . Observe, for any matrix  $M$  and projection matrix  $P$ ,  $M = PMP + (I-P)M(I-P) + PM(I-P) + (I-P)MP$ . Then,

$$\begin{aligned}
\|\Sigma_i - \hat{M}\|_F &\leq \underbrace{\|\hat{P}(\Sigma_i - \hat{M})\hat{P}\|_F}_{(1)} + \underbrace{\|(I - \hat{P})(\Sigma_i - \hat{M})(I - \hat{P})\|_F}_{(2)} \\
&\quad + \underbrace{\|\hat{P}(\Sigma_i - \hat{M})(I - \hat{P})\|_F}_{(3)} + \underbrace{\|(I - \hat{P})(\Sigma_i - \hat{M})\hat{P}\|_F}_{(4)}
\end{aligned} \tag{8.2}$$

We bound each of the terms above. Since  $\hat{P}(I + S - P(I + S)P)\hat{P} = 0$ , we can bound term (1) as follows

$$\|\hat{P}(\Sigma_i - \hat{M})\hat{P}\|_F = \|\hat{P}\Sigma_i\hat{P} - \hat{P}U\hat{\Sigma}_iU^\top\hat{P}\|_F = \|U^\top \Sigma_i U - \hat{\Sigma}_i\|_F \leq O\left(\frac{\sqrt{\Delta}\delta^{1/4}}{\alpha}\right) \tag{8.3}$$

Similarly, since  $(I - \hat{P})(U\hat{\Sigma}_iU^\top)(I - \hat{P}) = 0$  and  $\Sigma_i = I + \hat{S} - P_i - Q_i$ , we can bound term (2) as follows:

$$\begin{aligned}
\|(I - \hat{P})(\Sigma_i - \hat{M})(I - \hat{P})\|_F &= \|(I - \hat{P})(\Sigma_i - (I + \hat{S}))(I - \hat{P})\|_F \\
&\leq \|(I - \hat{P})(\Sigma_i - (I + \hat{S} - Q_i))(I - \hat{P})\|_F + \|(I - \hat{P})Q_i(I - \hat{P})\|_F \\
&\leq \|P_i\|_F^2 + \|(I - \hat{P})(Q_i - Q'_i)(I - \hat{P})\|_F + \|(I - \hat{P})Q'_i(I - \hat{P})\|_F \\
&\leq O\left(\sqrt{\frac{\eta}{\alpha}} + \frac{k^2\delta^{1/4}\Delta^{1/2}}{\alpha}\right)
\end{aligned} \tag{8.4}$$

Next, we bound term (3). Observe,  $\hat{P} \left( U \hat{\Sigma}_i U^\top \right) \left( I - \hat{P} \right) = 0$  and  $\hat{P} (I + S) \hat{P} (I - \hat{P}) = 0$ . Thus,

$$\begin{aligned}
\left\| \hat{P} \left( \Sigma_i - \hat{M} \right) \left( I - \hat{P} \right) \right\|_F &= \left\| \hat{P} \left( \Sigma_i - \left( I + \hat{S} \right) \right) \left( I - \hat{P} \right) \right\|_F \\
&= \left\| \hat{P} \left( P_i + Q_i \right) \left( I - \hat{P} \right) \right\|_F \\
&\leq \left\| \hat{P} P_i \left( I - \hat{P} \right) \right\|_F \\
&\leq O \left( \sqrt{\frac{\eta}{\alpha}} \right)
\end{aligned} \tag{8.5}$$

Observe, term (4) follows from a similar argument. Combining equations (8.3), (8.4), (8.5) and substituting back into (8.2) we can conclude

$$\left\| \Sigma_i - \hat{M} \right\|_F \leq O \left( \sqrt{\frac{\eta}{\alpha}} + \frac{k^2 \delta^{1/4} \Delta^{1/2}}{\alpha} \right)$$

The size of  $L'$  is  $\ell' = O \left( \log k \left( \eta / (k^5 (\Delta^4 + 1/\alpha^4)) \right)^{-4k} \right)$  and since we add a single tuple of  $k$  means and covariances for each subspace in  $L'$ , the list  $L$  has the same size. The running time is poly  $(|Y|, |L|, d^k, m, 1/\varepsilon_1)$  concluding the proof.  $\square$

## 8.2 Proof of Theorem 8.1

Since we have all the main ingredients: the tensor decomposition algorithm recovering a polynomial size of list (Theorem 8.2), the upgraded partial clustering algorithm with high probability of success (Theorem 7.1) and the spectral separation algorithm of thin components (Lemma 5.1), we can now complete the proof of Theorem 8.1.

The algorithm establishing Theorem 8.1 is almost the same as Algorithm 6.3. The only difference is we will replace Algorithm 4.4 by Algorithm 7.2 and replace Algorithm 3.2 by Algorithm 8.6. The following two lemmas show that by modifying the parameters slightly and applying the upgraded partial clustering and tensor decomposition algorithms, we can have the same conclusions as in Lemma 6.8 and Lemma 6.9 with a polynomial success probability. Then the proof of Theorem 8.1 is exactly the same as the proof of Theorem 6.1 in Section 6.2 except for the use of Lemma 8.10 and Lemma 8.11 instead of Lemma 6.8 and Lemma 6.9.

**Lemma 8.10** (Non-negligible Weight and Covariance Separation). *Given  $0 < \varepsilon < 1/k^{k^{O(k^2)}}$  and  $k \in \mathbb{N}$ , let  $\alpha = \varepsilon^{1/(45C^{k+1}(k+1)!)}$ .*

*Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with mixture covariance  $\Sigma$  such that  $w_i \geq \alpha$  for all  $i \in [k]$  and there exist  $i, j \in [k]$  such that  $\left\| \Sigma^{+1/2} (\Sigma_i - \Sigma_j) \Sigma^{+1/2} \right\|_F^2 > 1/\alpha^5$ . Further, let  $X$  be a set of points satisfying Condition 2.47 with respect to  $\mathcal{M}$  for some parameters  $\gamma \leq \varepsilon d^{-8k} k^{-Ck}$ , for a sufficiently large constant  $C$ , and  $t \geq 8k$ . Let  $Y$  be an  $\varepsilon$ -corrupted version of  $X$  of size  $n \geq n_0 = (dk)^{\Omega(1)} / \varepsilon$ , Algorithm 7.2 partitions  $Y$  into  $Y_1, Y_2$  in time  $n^{O(1)}$  such that with probability at least  $2^{-O(k)}(1 - O(\alpha))$  there is a non-trivial partition of  $[k]$  into  $Q_1 \cup Q_2$  so that letting  $\mathcal{M}_j$  be a distribution proportional to*



$\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then  $Y_j$  is an  $\mathcal{O}\left(\varepsilon^{1/(45C^{k+1}(k+1)!)}\right)$ -corrupted version of  $\bigcup_{i \in Q_j} X_i$  satisfying Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\mathcal{O}(k\gamma/W_j), t)$ .

*Proof.* We run Algorithm 7.2 with sample set  $Y$ , number of components  $k$ , the fraction of outliers  $\varepsilon$  and the accuracy parameter  $\eta$ . Since  $X$  satisfies Condition 2.47, we can set  $t = 10$ ,  $\beta = (k^2 t^4 \alpha)^{t/2} = O_k(\alpha^5)$  and  $\eta = \alpha^2 \gg \sqrt{\varepsilon/\alpha}$  in Theorem 7.1. Then, by assumption, there exist  $i, j$  such that

$$\|\Sigma^{t/2} (\Sigma_i - \Sigma_j) \Sigma^{t/2}\|_F^2 > \frac{1}{\alpha^5} = \Omega\left(\frac{k^2 t^4}{\beta^{2/t} \alpha^4}\right).$$

We observe that we also satisfy the other preconditions for Theorem 7.1, since  $n \geq (dk/\varepsilon)^{\Omega(1)}$ .

Then, Theorem 7.1 implies that with probability at least  $2^{-O(k)}(1 - O(\eta/\alpha - \sqrt{\eta})) = 2^{-O(k)}(1 - O(\alpha))$ , the set  $Y$  is partitioned in two sets  $Y_1$  and  $Y_2$  such that there is a non-trivial partition of  $[k]$  into  $Q_1 \cup Q_2$  so that letting  $\mathcal{M}_j$  be a distribution proportional to  $\sum_{i \in Q_j} w_i G_i$  and  $W_j = \sum_{i \in Q_j} w_i$ , then  $Y_j$  is an  $\mathcal{O}\left(\varepsilon^{1/(45C^{k+1}(k+1)!)}\right)$ -corrupted version of  $\bigcup_{i \in Q_j} X_i$ . By Lemma 2.50,  $\bigcup_{i \in Q_j} X_i$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\mathcal{O}(k\gamma/W_j), t)$ .  $\square$

When the mixture is not covariance separated and nearly isotropic, we can obtain a small list of hypotheses such that one of them is close to the true parameters, via tensor decomposition.

**Lemma 8.11** (Mixture is List-decodable). *Given  $0 < \varepsilon < 1/k^{k\mathcal{O}(k^2)}$  let  $\alpha = \varepsilon^{1/(45C^{k+1}(k+1)!)}$ . Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  with  $G_i = \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with mixture mean  $\mu$  and mixture covariance  $\Sigma$ , such that  $\|\mu\|_2 \leq \mathcal{O}(\sqrt{\varepsilon/\alpha})$ ,  $\|\Sigma - I\|_F \leq \mathcal{O}(\sqrt{\varepsilon/\alpha})$ ,  $w_i \geq \alpha$  for all  $i \in [k]$ , and  $\|\Sigma_i - \Sigma_j\|_F^2 \leq 1/\alpha^5$  for any pair of components, and let  $X$  be a set of points satisfying Condition 2.47 with respect to  $\mathcal{M}$  for some parameters  $\gamma = \varepsilon d^{-8k} k^{-Ck}$ , for a sufficiently large constant  $C$ , and  $t = 8k$ . Let  $Y$  be an  $\varepsilon$ -corrupted version of  $X$  of size  $n$ , Algorithm 8.6 outputs a list  $L$  of hypotheses of size  $O((1/\varepsilon)^{4k^2})$  in time  $\text{poly}(|L|, n)$  such that if we choose a hypothesis  $\{\hat{\mu}_i, \hat{\Sigma}_i\}_{i \in [k]}$  uniformly at random,  $\|\mu_i - \hat{\mu}_i\|_2 \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$  and  $\|\Sigma_i - \hat{\Sigma}_i\|_F \leq \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right)$  for all  $i$  with probability at least  $O(\varepsilon^{4k^2})$ .*

*Proof.* Recall we run Algorithm 8.6 on the samples  $Y$ , the number of clusters  $k$ , the fraction of outliers  $\varepsilon$  and the minimum weight  $\alpha = \varepsilon^{1/(20C^{k+1}(k+1)!)}$ . Next, we show that the preconditions of Theorem 8.2 are satisfied. First, the upper bounds on  $\|\mu\|_2$  and  $\|\Sigma - I\|_F$  imply  $\sum_{i \in [k]} w_i (\Sigma_i + \mu_i \mu_i^\top) = \Sigma + \mu \mu^\top \leq (1 + \mathcal{O}(\sqrt{\varepsilon/\alpha}))I$ . Since the LHS is a conic combination of PSD matrices, it follows that for all  $i \in [k]$ ,  $\mu_i \mu_i^\top \leq \frac{1}{\alpha} (1 + \mathcal{O}(\sqrt{\varepsilon/\alpha}))I$ , and thus  $\|\mu_i \mu_i^\top\|_F \leq \frac{2}{\alpha}$ . Next, we can write:

$$\begin{aligned} \|\Sigma_i - I\|_F &\leq \|\Sigma_i - (\Sigma + \mu \mu^\top)\|_F + \|\Sigma - I\|_F + \|\mu \mu^\top\|_F \\ &= \left\| \Sigma_i - \sum_{j \in [k]} w_j (\Sigma_j + \mu_j \mu_j^\top) \right\|_F + \frac{\sqrt{\varepsilon}k}{\alpha} + \frac{\varepsilon}{\alpha} \\ &\leq \left\| \sum_{j \in [k]} w_j (\Sigma_i - \Sigma_j) \right\|_F + \frac{2}{\alpha} + \frac{\sqrt{\varepsilon}k}{\alpha} + \frac{\varepsilon}{\alpha} \end{aligned}$$

$$\leq \frac{2}{\alpha^{5/2}},$$

where the first and the third inequalities follow from the triangle inequality and the upper bound on  $\|\mu_i \mu_i^\top\|_F$ , and the last inequality follows from the assumption that  $\|\Sigma_i - \Sigma_j\|_F^2 \leq 1/\alpha^5$  for every pair of covariances  $\Sigma_i, \Sigma_j$ . So, we can set  $\Delta = 2\alpha^{-5/2}$  in Theorem 8.2. Then, given the definition of  $\alpha$ , we have that

$$\eta = 2k^{4k} \mathcal{O}(1 + \Delta/\alpha)^{4k} \sqrt{\varepsilon} = \mathcal{O}(\varepsilon^{2/5})$$

and  $1/\varepsilon^2 \geq \log(1/\eta)(k + 1/\alpha + \Delta)^{4k+5}/\eta^2$ . Therefore, Algorithm 8.6 outputs a list  $L$  of hypotheses such that  $|L| = \exp(1/\varepsilon^2)$ , and with probability at least 0.99,  $L$  contains a hypothesis that satisfies the following: for all  $i \in [k]$ ,

$$\begin{aligned} \|\hat{\mu}_i - \mu_i\|_2 &= \mathcal{O}\left(\frac{\Delta^{1/2}}{\alpha}\right) \eta^{G(k)} = \mathcal{O}\left(\varepsilon^{-1/(20C^{k+1}(k+1)!)} \cdot \varepsilon^{1/(10C^{k+1}(k+1)!)}\right) = \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right) \text{ and} \\ \|\hat{\Sigma}_i - \Sigma_i\|_F &= \mathcal{O}(k^4) \frac{\Delta^{1/2}}{\alpha} \eta^{G(k)} = \mathcal{O}\left(\varepsilon^{1/(20C^{k+1}(k+1)!)}\right). \end{aligned} \quad (8.6)$$

Then if we choose a hypothesis in  $L$  uniformly at random, the probability that we choose the hypothesis satisfying (6.2) is at least  $1/|L| = \exp(-1/\varepsilon^2)$ .  $\square$

## 9 Robust Parameter Recovery: Proof of Theorem 1.6

In order to show that our algorithm recovers the individual components and the parameters, we will prove the following identifiability theorem. Without any assumption on the mixtures, it is impossible to distinguish components within  $\varepsilon$  total variation distance with  $\varepsilon$ -fraction of noise. So given two mixtures of Gaussians with  $\varepsilon$  total variation distance, the theorem shows that there exist two partitions of components of the two mixtures respectively such that any two components in the matched pair are  $\text{poly}(\varepsilon)$ -close in total variation distance.

**Theorem 9.1** (Identifiability). *Let  $\mathcal{M} = \sum_{i=1}^{k_1} w_i G_i$ ,  $\mathcal{M}' = \sum_{i=1}^{k_2} w'_i G'_i$  be two mixtures of Gaussians such that  $d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \leq \varepsilon$ . Then there exists a partition of  $[k_1]$  into sets  $R_0, R_1, \dots, R_\ell$  and a partition of  $[k_2]$  into sets  $S_0, S_1, \dots, S_\ell$  such that*

1. Let  $W_i = \sum_{j \in R_i} w_j$  for  $i = 0, 1, \dots, k_1$ ,  $W'_i = \sum_{j \in S_i} w'_j$  for  $i = 0, 1, \dots, k_2$ . Then for all  $i \in [\ell]$ ,

$$\begin{aligned} |W_i - W'_i| &\leq \text{poly}_k(\varepsilon) \\ d_{\text{TV}}(G_j, G'_{j'}) &\leq \text{poly}_k(\varepsilon) \quad \forall j \in R_i, j' \in S_i \end{aligned}$$

2.  $W_0, W'_0 \leq \text{poly}_k(\varepsilon)$ .

**Corollary 9.2.** *There is an algorithm with the following behavior: Given  $\varepsilon > 0$  and a multiset of  $n = d^{O(k)} \text{poly}(\varepsilon)$  samples from a distribution  $F$  on  $\mathbb{R}^d$  such that  $d_{\text{TV}}(F, \mathcal{M}) \leq \varepsilon$ , for an unknown target  $k$ -GMM  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$ , the algorithm runs in time  $d^{O(k)} \text{poly}_k(1/\varepsilon)$  and outputs a  $k'$ -GMM hypothesis  $\widehat{\mathcal{M}} = \sum_{i=1}^{k'} \widehat{w}_i \mathcal{N}(\widehat{\mu}_i, \widehat{\Sigma}_i)$  with  $k' \leq k$  such that with high probability there exists a partition of  $[k]$  into  $k' + 1$  sets  $R_0, R_1, \dots, R_{k'}$  such that*

1. Let  $W_i = \sum_{j \in R_i} w_j$ . Then for all  $i \in [k']$ ,

$$|W_i - \hat{w}_i| \leq \text{poly}_k(\varepsilon)$$

$$d_{\text{TV}}(\mathcal{N}(\mu_j, \Sigma_j), \mathcal{N}(\hat{\mu}_i, \hat{\Sigma}_i)) \leq \text{poly}_k(\varepsilon) \quad \forall j \in R_i$$

2. The sum of weights of exceptional components in  $R_0$  is at most  $\text{poly}_k(\varepsilon)$ .

Parameter estimation is implied by TV distance for individual Gaussians (in relative Frobenius norm). The corollary follows immediately from the identifiability theorem.

**Outline of Proof.** The first step is to deal with the components in  $\mathcal{M}$  and  $\mathcal{M}'$  with small weights. We will construct  $\tilde{\mathcal{M}}, \tilde{\mathcal{M}'}$  by removing components with small weights. If we prove the statement on  $\tilde{\mathcal{M}}, \tilde{\mathcal{M}'}$ , we can then deduce the theorem in the general case with worse, but still polynomial dependencies on  $\varepsilon$ . The second step is a partial clustering, after which the components within each cluster have TV distance bounded by  $1 - \text{poly}(\varepsilon)$ . We prove this lemma in a separate section. After that we modify the parameters slightly so that the resulting parameters for different components are either identical or have a minimum separation. After this, we can use a lemma from [LM20] that provides a 1-1 mapping between the components of two such mixtures with small TV distance such that the mapped pairs have small TV distance.

**Distance between Gaussians.** We use the following facts for Gaussian distributions.

**Lemma 9.3** (Frobenius Distance to TV Distance). *Suppose  $N(\mu_1, \Sigma_1), N(\mu_2, \Sigma_2)$  are Gaussians with  $\|\mu_1 - \mu_2\|_2 \leq \delta$  and  $\|\Sigma_1 - \Sigma_2\|_F \leq \delta$ . If the eigenvalues of  $\Sigma_1$  and  $\Sigma_2$  are at least  $\lambda > 0$ , then  $d_{\text{TV}}(N(\mu_1, \Sigma_1), N(\mu_2, \Sigma_2)) = O(\delta/\lambda)$ .*

**Lemma 9.4** (Lemma 5.4 in [LM20]). *Let  $\mathcal{M}$  be a mixture of  $k$  Gaussians that is connected if we draw edges between all components  $i, j$  in  $\mathcal{M}$  such that  $d_{\text{TV}}(G_i, G_j) \leq 1 - \delta$ . Let  $\Sigma$  be the covariance matrix of  $\mathcal{M}$ . Then for any components  $\Sigma_i$  of the mixture*

1.  $\Sigma_i \geq \text{poly}_k(\delta)\Sigma$
2.  $\|\Sigma^{-1/2}(\Sigma - \Sigma_i)\Sigma^{-1/2}\|_F \leq \text{poly}_k(\delta)^{-1}$ .

The proof is identical to Lemma 5.4 in [LM20]. The only difference is that in [LM20] the authors assume that the minimal weight of  $\mathcal{M}$  is at least  $\delta$  and TV distance between any pair of components is at least  $\delta$  but here we do not need these two assumptions, which does not affect the proof.

**Fact 9.5** (Claim 3.9 in [LM20]). *Let  $\partial$  denote the differential operator with respect to  $y$ . If*

$$f(y) = P(y, X) \exp\left(a(X)y + \frac{1}{2}b(X)y^2\right)$$

where  $P$  is a polynomial in  $y$  of degree  $k$  (whose coefficients are polynomials in  $X$ ) and  $a(X), b(X)$  are polynomials in  $X$  then

$$(\partial - (a(X) + yb(X)))f(y) = Q(y, X) \exp\left(a(X)y + \frac{1}{2}b(X)y^2\right)$$

where  $Q$  is a polynomial in  $y$  with degree exactly  $k - 1$  whose leading coefficient is  $k$  times the leading coefficient of  $P$ .

**Fact 9.6** (Corollary 3.10 in [LM20]). Let  $\partial$  denote the differential operator with respect to  $y$ . If

$$f(y) = P(y, X) \exp\left(a(X)y + \frac{1}{2}b(X)y^2\right)$$

where  $P$  is a polynomial in  $y$  of degree  $k$  then

$$(\partial - (a(X) + yb(X)))^{k+1} f(y) = 0.$$

**Fact 9.7** (Claim 3.11 in [LM20]). Let  $\partial$  denote the differential operator with respect to  $y$ . If

$$f(y) = P(y, X) \exp\left(a(X)y + \frac{1}{2}b(X)y^2\right)$$

where  $P$  is a polynomial in  $y$  of degree  $k$ . Let the leading coefficient of  $P$  (viewed as a polynomial in  $y$ ) be  $L(X)$ . Let  $c(X)$  be a linear polynomial in  $X$  and  $d(X)$  be a quadratic polynomial in  $X$  such that  $\{a(X), b(X)\} \neq \{c(X), d(X)\}$ . If  $b(X) \neq d(X)$  then

$$(\partial - (c(X) + yd(X)))^{k'} f(y) = Q(y, X) \exp\left(a(X)y + \frac{1}{2}b(X)y^2\right)$$

where  $Q$  is a polynomial of degree  $k + k'$  in  $y$  with leading coefficient

$$L(X)(b(X) - d(X))^{k'}$$

and if  $b(X) = d(X)$  then

$$(\partial - (c(X) + yd(X)))^{k'} f(y) = Q(y, X) \exp\left(a(X)y + \frac{1}{2}b(X)y^2\right)$$

where  $Q$  is a polynomial of degree  $k$  in  $y$  with leading coefficient

$$L(X)(a(X) - c(X))^{k'}.$$

**Lemma 9.8.** Let  $\mathcal{M} = \sum_{i=1}^{k_1} w_i G_i$ ,  $\mathcal{M}' = \sum_{i=1}^{k_2} w'_i G'_i$  be two mixtures of Gaussians such that  $d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \leq \varepsilon$ . For any constant  $0 < c_1 < 1$ , there exists  $i \in [k_1 + k_2 + 1]$  such that  $w_j, w'_j \notin [\varepsilon^{c_1^{i-1}}, \varepsilon^{c_1^i}]$  for any  $j \in [k_1], j' \in [k_2]$ . Moreover, if

$$\begin{aligned} \tilde{\mathcal{M}} &= \frac{\sum_{\{j:w_j \geq \varepsilon^{c_1^i}\}} w_j G_j}{\sum_{\{j:w_j \geq \varepsilon^{c_1^i}\}} w_j} \\ \tilde{\mathcal{M}}' &= \frac{\sum_{\{j:w'_j \geq \varepsilon^{c_1^i}\}} w'_j G'_j}{\sum_{\{j:w'_j \geq \varepsilon^{c_1^i}\}} w'_j} \end{aligned}$$

then  $d_{\text{TV}}(\tilde{\mathcal{M}}, \tilde{\mathcal{M}}') \leq O_k(\varepsilon^{c_1^{i-1}})$ .

*Proof.* We can see that  $[\varepsilon^{c_1^{i-1}}, \varepsilon^{c_1^i})$  with  $i \in [k_1 + k_2 + 1]$  are  $k_1 + k_2 + 1$  disjoint intervals and  $w_j, w'_j$ , with  $j \in [k_1], j' \in [k_2]$  have at most  $k_1 + k_2$  distinct values. So there is one interval containing no weights.

We then construct  $\tilde{\mathcal{M}}$  by removing the small components in  $\mathcal{M}$ . The sum of weights removed is at most  $k\varepsilon^{c_1^{i-1}}$ . So  $d_{\text{TV}}(\mathcal{M}, \tilde{\mathcal{M}}) \leq k\varepsilon^{c_1^{i-1}}$ . Similarly, we have  $d_{\text{TV}}(\mathcal{M}', \tilde{\mathcal{M}}') \leq k\varepsilon^{c_1^{i-1}}$ . By the triangle inequality,

$$d_{\text{TV}}(\tilde{\mathcal{M}}, \tilde{\mathcal{M}}') \leq d_{\text{TV}}(\mathcal{M}, \mathcal{M}') + d_{\text{TV}}(\mathcal{M}, \tilde{\mathcal{M}}) + d_{\text{TV}}(\mathcal{M}', \tilde{\mathcal{M}}') \leq O_k(\varepsilon^{c_1^{i-1}}).$$

□

Lemma 9.8 shows that we can remove components with tiny weights in the mixtures. So in the following lemma, we will assume  $\mathcal{M}$  and  $\mathcal{M}'$  are Gaussian mixtures with minimal weights at least  $\text{poly}(\varepsilon)$ . We will show that we can partition the union of components of two mixtures so that if we prove Theorem 9.1 for each part of the partition, we can combine them to prove Theorem 9.1 on the full mixtures.

**Lemma 9.9.** *For any constant  $0 < c_3 < 1$ , there exist  $c_1, c_2 > 0$  that depend on  $k$  and  $c_3$ , such that if  $\mathcal{M} = \sum_{i=1}^{k_1} w_i G_i, \mathcal{M}' = \sum_{i=1}^{k_2} w'_i G'_i$  with  $k_1, k_2 \leq k, d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \leq \varepsilon$  and  $w_i, w'_i \geq \varepsilon^{c_1}$  for all  $i$ , then there exists a partition of  $[k_1]$  into sets  $R_1, \dots, R_\ell$  and a partition of  $[k_2]$  into sets  $S_1, \dots, S_\ell$  such that*

1. *For all  $i \in [\ell]$ , let  $W_i = \sum_{j \in R_i} w_j, W'_i = \sum_{j \in S_i} w'_j$  be the sum of weights in each piece. Let  $\mathcal{M}_i = \frac{1}{W_i} \sum_{j \in R_i} w_j G_j, \mathcal{M}'_i = \frac{1}{W'_i} \sum_{j \in S_i} w'_j G'_j$  be the submixtures of Gaussians after partition. Then for all  $i \in [\ell]$ ,*

$$\begin{aligned} |W_i - W'_i| &\leq \text{poly}_k(\varepsilon) \\ d_{\text{TV}}(\mathcal{M}_i, \mathcal{M}'_i) &\leq O_k(\varepsilon^{c_2}) \end{aligned}$$

2. *Consider the graph with vertices corresponding to components in  $\mathcal{M}$  and  $\mathcal{M}'$  and two components are adjacent if the total variation distance between them is at most  $1 - \varepsilon^{c_2 c_3}$ . Then the induced subgraph of vertices with indices  $R_i \cup S_i$  is connected for all  $i \in [\ell]$ .*

The proof of Lemma 9.9 is deferred to Section 9.1. In the following two lemmas, we then prove Theorem 9.1 for each pair  $\mathcal{M}_i, \mathcal{M}'_i$  defined in Lemma 9.9. In Lemma 9.10, we construct two mixtures of which pairs of parameters are identical or separated. We also shows it suffices to work under this simplification.

**Lemma 9.10.** *For any constant  $0 < c_4 < 1$ , there exist  $c_3, c_5$  that depend on  $k$  and  $c_4$ , such that if  $\mathcal{M} = \sum_{i=1}^{k_1} w_i G_i, \mathcal{M}' = \sum_{i=1}^{k_2} w'_i G'_i$  with  $k_1, k_2 \leq k$  and*

1.  $\frac{1}{2}\mathcal{M} + \frac{1}{2}\mathcal{M}'$  is isotropic,
2.  $d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \leq \varepsilon$ ,
3.  $w_i, w'_i \geq \varepsilon^{c_3}$  for all  $i$ ,

4. Let  $\mathcal{G}$  be a graph with components  $G_i, G'_i$  in  $\mathcal{M}$  and  $\mathcal{M}'$  as vertex set and two components are adjacent if the total variation distance between them is at most  $1 - \varepsilon^{c_3}$ . Then  $\mathcal{G}$  is connected

then there exist two mixtures of Gaussians  $\tilde{\mathcal{M}} = \sum_{i=1}^{\tilde{k}_1} \tilde{w}_i \tilde{G}_i$ ,  $\tilde{\mathcal{M}}' = \sum_{i=1}^{\tilde{k}_2} \tilde{w}'_i \tilde{G}'_i$  such that

1. Any pair in  $\{\tilde{\mu}_i\} \cup \{\tilde{\mu}'_i\}$  is either identical or separated by at least  $\varepsilon^{c_4 c_5}$
2. Any pair in  $\{\tilde{\Sigma}_i\} \cup \{\tilde{\Sigma}'_i\}$  is either identical or separated by at least  $\varepsilon^{c_4 c_5}$  in Frobenius norm.
3.  $\left\| \mathbb{E}(h_m(\tilde{\mathcal{M}})) - \mathbb{E}(h_m(\tilde{\mathcal{M}}')) \right\|_F \leq O_k(\varepsilon^{c_5})$  for any  $m \leq O(k)$
4. There exist  $\pi_1 : [k_1] \rightarrow [\tilde{k}_1]$  and  $\pi_2 : [k_2] \rightarrow [\tilde{k}_2]$  such that

$$\begin{aligned} \sum_{i:\pi_1(i)=j} w_i &= \tilde{w}_j, & \sum_{i:\pi_2(i)=j} w'_i &= \tilde{w}'_j, \\ d_{\text{TV}}(G_i, \tilde{G}_{\pi_1(i)}) &\leq \text{poly}_k(\varepsilon), & \text{for all } i \in [k_1] \\ d_{\text{TV}}(G'_i, \tilde{G}'_{\pi_2(i)}) &\leq \text{poly}_k(\varepsilon), & \text{for all } i \in [k_2]. \end{aligned}$$

*Proof.* For any  $0 < c_4 < 1$ , there is  $\ell \in [k^2]$  such that the distance between any pair of parameters in  $\{\mu_i\} \cup \{\mu'_i\}$  or the Frobenius distance between any pair in  $\{\Sigma_i\} \cup \{\Sigma'_i\}$  is not in the interval  $[\varepsilon^{(c_4/2)^{\ell-1}}, \varepsilon^{(c_4/2)^\ell}]$ .

Now consider a graph  $\mathcal{G}$  on  $k_1 + k_2$  nodes where each node represents a vector in  $\{\mu_i\} \cup \{\mu'_i\}$  and two vectors  $a, b$  are adjacent if

$$\|a - b\| \leq \varepsilon^{(c_4/2)^{\ell-1}}.$$

We now construct new mixtures  $\tilde{\mathcal{M}}, \tilde{\mathcal{M}}'$ . For each connected component in  $\mathcal{G}$  say  $\{\mu_{i_1}, \dots, \mu'_{j_1}, \dots\}$ , pick a representative say  $\mu_{i_1}$  and set  $\tilde{\mu}_{i_1} = \dots = \tilde{\mu}'_{j_1} = \dots = \mu_{i_1}$ . Do this for all connected components and similar in the graph on covariance matrices with edges  $(i, j)$  if

$$\|\Sigma_i - \Sigma_j\|_F \leq \varepsilon^{(c_4/2)^{\ell-1}}.$$

After replacing close parameters with a representative, we may get some exactly same components in each new mixture. We then merge components with same means and covariances by adding their weights. Since all representatives of means and covariances are in different connected components of the graphs, they are separated by at least  $\varepsilon^{(c_4/2)^\ell}$ . Setting  $c_5 = 1/2(c_4/2)^{\ell-1}$  gives a separation of  $\varepsilon^{c_4 c_5}$ .

Next we prove 3. There is a natural mapping  $\pi_1 : [k_1] \rightarrow [\tilde{k}_1]$  that maps any component in  $\mathcal{M}$  to the merged component in  $\tilde{\mathcal{M}}$  and a similar mapping  $\pi_2 : [k_2] \rightarrow [\tilde{k}_2]$  for  $\mathcal{M}', \tilde{\mathcal{M}}'$ . For all  $i$ , we have

$$\left\| \tilde{\mu}_{\pi_1(i)} - \mu_i \right\|, \left\| \tilde{\mu}'_{\pi_2(i)} - \mu'_i \right\|, \left\| \tilde{\Sigma}_{\pi_1(i)} - \Sigma_i \right\|_F, \left\| \tilde{\Sigma}'_{\pi_2(i)} - \Sigma'_i \right\|_F \leq O_k(1) \varepsilon^{(c_4/2)^{\ell-1}} \quad (9.1)$$

because for any pair of parameters above say  $\tilde{\mu}_{\pi_1(i)}$  and  $\mu_i$ , there is a path of length at most  $2k$  connecting  $\mu_i$  to the representative of the connected component, and each edge connects a pair

with TV distance at most  $\varepsilon$ . Suppose  $\|\mu_i\|, \|\Sigma_i - I\|_F \leq \Delta$ . Then by Definition 2.4, we have for any integer  $m$ ,

$$\left\| \mathbb{E}(h_m(\mathcal{M})) - \mathbb{E}(h_m(\tilde{\mathcal{M}})) \right\|_F \leq O_k(m) \Delta^m \varepsilon^{(c_4/2)^{\ell-1}}.$$

Since  $\frac{1}{2}\mathcal{M} + \frac{1}{2}\mathcal{M}'$  is isotropic and the minimum weight in  $\frac{1}{2}\mathcal{M} + \frac{1}{2}\mathcal{M}'$  is at least  $\frac{1}{2}\varepsilon^{c_3}$ , we have  $\|\mu_i\| \leq \sqrt{2/\varepsilon^{c_3}}$  for all  $i$ . Applying Lemma 9.4 to  $\frac{1}{2}\mathcal{M} + \frac{1}{2}\mathcal{M}'$ , we have  $\|I - \Sigma_i\|_F \leq \text{poly}_k(\varepsilon^{c_3})^{-1}$ . So there is a constant  $a$  such that  $\Delta \leq \varepsilon^{-ac_3}$ . If we take  $c_3 > 0$  so that  $ac_3 O(k) \leq 1/2(c_4/2)^{\ell-1}$  and take  $c_5 = 1/2(c_4/2)^{\ell-1}$ , then

$$\left\| \mathbb{E}(h_m(\mathcal{M})) - \mathbb{E}(h_m(\tilde{\mathcal{M}})) \right\|_F \leq O_k(m) \varepsilon^{(c_4/2)^{\ell-1} - O(m)ac_3} = O_k(\varepsilon^{c_5})$$

for  $m \leq O(k)$ . By the same argument, we have the similar inequality for  $\mathcal{M}'$  and  $\tilde{\mathcal{M}}'$

$$\left\| \mathbb{E}(h_m(\mathcal{M}')) - \mathbb{E}(h_m(\tilde{\mathcal{M}}')) \right\|_F = O_k(\varepsilon^{c_5}).$$

Since we can use Proposition 3.3 to robustly estimate the Hermite tensors of a Gaussian mixture with  $\varepsilon$ -fraction of noise and  $\text{poly}(\varepsilon)$  error guarantee, we must have

$$\|\mathbb{E}(h_m(\mathcal{M})) - \mathbb{E}(h_m(\mathcal{M}'))\|_F \leq \text{poly}_k(\varepsilon).$$

Then by the triangle inequality,

$$\begin{aligned} \left\| \mathbb{E}(h_m(\tilde{\mathcal{M}})) - \mathbb{E}(h_m(\tilde{\mathcal{M}}')) \right\|_F &\leq \left\| \mathbb{E}(h_m(\mathcal{M})) - \mathbb{E}(h_m(\tilde{\mathcal{M}})) \right\|_F + \\ &\quad \|\mathbb{E}(h_m(\mathcal{M})) - \mathbb{E}(h_m(\mathcal{M}'))\|_F + \left\| \mathbb{E}(h_m(\mathcal{M}')) - \mathbb{E}(h_m(\tilde{\mathcal{M}}')) \right\|_F = O(\varepsilon^{c_5}). \end{aligned}$$

For the last conclusion, from the definition of  $\pi_1$  and  $\pi_2$ , we know that

$$\sum_{i:\pi_1(i)=j} w_i = \tilde{w}_j, \quad \sum_{i:\pi_2(i)=j} w'_i = \tilde{w}'_j.$$

Applying Lemma 9.4 to  $\frac{1}{2}\mathcal{M} + \frac{1}{2}\mathcal{M}'$ , we have that eigenvalues of  $\Sigma_i$  and  $\Sigma'_i$  are at least  $\text{poly}(\varepsilon^{c_3})$  for all  $i$ . Then if  $c_3$  is sufficiently small, by Lemma 9.3, (9.1) implies  $d_{\text{TV}}(G_i, \tilde{G}_{\pi(i)}) \leq \text{poly}_k(\varepsilon)$  and  $d_{\text{TV}}(G'_i, \tilde{G}'_{\pi(i)}) \leq \text{poly}_k(\varepsilon)$  for all  $i$ .  $\square$

The following lemma shows the identifiability under the simplification of Lemma 9.10. It is proved in the proof of Lemma 8.2 in [LM20].

**Lemma 9.11.** *Suppose  $\mathcal{M} = \sum_{i=1}^{k_1} w_i G_i$ ,  $\mathcal{M}' = \sum_{i=1}^{k_2} w'_i G'_i$  satisfies 1,2,3 in the conclusion of Lemma 9.10 with constants  $c_4, c_5$  and the minimal weights are at least  $\varepsilon^{c_3}$ . There exists a sufficiently small function  $f(k) > 0$  depending only on  $k$  such that if  $c_4 \leq f(k)$ , then  $k_1 = k_2$  and there exists a permutation  $\pi$  such that  $|w_i - w'_{\pi(i)}| \leq \text{poly}_k(\varepsilon)$  and  $G_i = G'_{\pi(i)}$ .*



*Proof.* Consider the component  $G'_{k_2} = N(\mu'_{k_2}, \Sigma'_{k_2})$  in  $\mathcal{M}'$ . We claim that there must be some  $i \in [k_1]$  such that

$$(\mu_i, \Sigma_i) = (\mu'_{k_2}, \Sigma'_{k_2}).$$

Assume for the sake of contradiction that this is not the case. Let  $S_1 = \{i \in [k_1] : \Sigma_i = \Sigma'_{k_2}\}$  and  $S_2 = \{i \in [k_2 - 1] : \Sigma'_i = \Sigma'_{k_2}\}$ . Suppose  $F, F'$  are the generating functions of  $\mathcal{M}$  and  $\mathcal{M}'$

$$F = \sum_{i=1}^{k_1} w_i \exp\left(\mu_i^T X + \frac{1}{2} X^T \Sigma_i X y^2\right) = \sum_{m=0}^{\infty} \frac{1}{m!} h_m(\mathcal{M}) y^m$$

$$F' = \sum_{i=1}^{k_2} w'_i \exp\left(\mu'_i{}^T X + \frac{1}{2} X^T \Sigma'_i X y^2\right) = \sum_{m=0}^{\infty} \frac{1}{m!} h_m(\mathcal{M}') y^m.$$

Then define the differential operators

$$\mathcal{D}_i = \partial - \mu_i^T X - X^T \Sigma_i X y$$

$$\mathcal{D}'_i = \partial - \mu'_i{}^T X - X^T \Sigma'_i X y$$

where partial derivatives are taken with respect to  $y$ . Now consider the differential operator

$$\mathcal{D} = (\mathcal{D}'_{k_2-1})^{2^{k_1+k_2-2}} \dots (\mathcal{D}'_1)^{2^{k_1}} \mathcal{D}_{k_1}^{2^{k_1-1}} \dots \mathcal{D}_1$$

By Fact 9.6,  $\mathcal{D}(F) = 0$ . By Fact 9.6 and Fact 9.7, we have

$$\mathcal{D}(F') = P(y, X) \exp\left(\mu'_{k_2}{}^T X + \frac{1}{2} X^T \Sigma'_{k_2} X y^2\right)$$

where  $P$  is a polynomial of degree

$$\deg(P) = 2^{k_1+k_2-1} - 1 - \sum_{i \in S_1} 2^{i-1} - \sum_{i \in S_2} 2^{k_1+i-2}$$

with leading coefficient

$$C_0 = w'_{k_2} \prod_{i \in [k_1] \setminus S_1} (X^T (\Sigma'_{k_2} - \Sigma_i) X)^{2^{i-1}} \prod_{i \in S_1} ((\mu'_{k_2} - \mu_i)^T X)^{2^{i-1}}$$

$$\prod_{i \in [k_2-1] \setminus S_2} (X^T (\Sigma'_{k_2} - \Sigma'_i) X)^{2^{k_1+i-2}} \prod_{i \in S_2} ((\mu'_{k_2} - \mu'_i)^T X)^{2^{k_1+i-2}}.$$

We now compare the following differentials evaluated at  $y = 0$

$$(\mathcal{D}'_{k_2})^{\deg(P)} \mathcal{D}(F)$$

$$(\mathcal{D}'_{k_2})^{\deg(P)} \mathcal{D}(F')$$

The first quantity is 0 because  $\mathcal{D}(F)$  is identically 0 as a formal power series. The second one is  $\Omega_k(1)C_0$ . Since for any  $i$   $(\mu_i, \Sigma_i) \neq (\mu'_{k_2}, \Sigma'_{k_2})$ , our assumptions imply that the separation between

$\mu_i, \mu'_{k_2}$  or  $\Sigma_i, \Sigma'_{k_2}$  is at least  $\varepsilon^{c_4 c_5}$ . Then we have  $C_0 \geq \varepsilon^{c_4 c_5 O_k(1)}$  for some  $X$ . On the other hand, the coefficients of the formal power series  $F, F'$  are the Hermite polynomials  $h_m(\mathcal{M})$  and  $h_m(\mathcal{M}')$ . This is a contradiction with our assumption that

$$\|\mathbb{E}(h_m(\mathcal{M}) - \mathbb{E}(h_m(\mathcal{M}))\|_F \leq O_k(\varepsilon^{c_5})$$

as long as  $c_4$  is smaller than some sufficiently small function  $f(k)$  depending only on  $k$ . Thus there must be some component of  $\mathcal{M}$  that matches  $G'_{k_2} = N(\mu'_{k_2}, \Sigma'_{k_2})$ . We can repeat the argument for each component in  $\mathcal{M}'$  and in  $\mathcal{M}$  to conclude that  $\mathcal{M}$  and  $\mathcal{M}'$  have the same components.

Next we will show that the weights of the same components in  $\mathcal{M}$  and  $\mathcal{M}'$  are close. We can assume that  $\mathcal{M} = \sum_{i=1}^k w_i G_i, \mathcal{M}' = \sum_{i=1}^k w'_i G_i$  are two mixtures on the same set of components. Without loss of generality,

$$w_1 - w'_1 \leq \dots \leq w_\ell - w'_\ell \leq 0 \leq w_{\ell+1} - w'_{\ell+1} \leq \dots \leq w_k - w'_k.$$

Then we can consider the following two mixtures

$$\begin{aligned} & (w_1 - w'_1)G_1 + \dots + (w_\ell - w'_\ell)G_\ell \\ & (w_{\ell+1} - w'_{\ell+1})G_{\ell+1} + \dots + (w_k - w'_k)G_k. \end{aligned}$$

If

$$\sum_{i=1}^k |w_i - w'_i| > \varepsilon^\zeta$$

for some sufficiently small  $\zeta$  depending only on  $k$ , we can then normalize each of the above into a distribution and repeat the same argument, using the fact that pairs of components cannot be too close, to obtain a contradiction. Thus, the mixing weights of  $\mathcal{M}$  and  $\mathcal{M}'$  are  $\text{poly}_k(\varepsilon)$ -close and this completes the proof.  $\square$

*Proof of Theorem 9.1.* We first set  $c_4 = f(k)$  as in Lemma 9.11, and then  $c_3, c_5$  according to  $c_4$  as in Lemma 9.10, and  $c'_1, c_2$  according to  $c_3$  as in Lemma 9.9. Let  $c_1 = \min\{c'_1, c_2 c_3\}$ .

By Lemma 9.8, we can find  $i$  such that there is no  $w_j, w'_j$  in  $[\varepsilon^{c_1^{i-1}}, \varepsilon^{c_1^i}]$ . Let  $\tilde{\mathcal{M}} = \sum_{\{j:w_j \geq \varepsilon^{c_1^i}\}} w_j G_j$  and  $\tilde{\mathcal{M}}' = \sum_{\{j:w'_j \geq \varepsilon^{c_1^i}\}} w'_j G'_j$ . Then  $d_{\text{TV}}(\tilde{\mathcal{M}}, \tilde{\mathcal{M}}') \leq O(\varepsilon^{c_1^{i-1}})$ . Let  $\varepsilon_1 = \varepsilon^{c_1^{i-1}}$ . We have  $d_{\text{TV}}(\tilde{\mathcal{M}}, \tilde{\mathcal{M}}') \leq O(\varepsilon_1)$  and the minimum weights of  $\tilde{\mathcal{M}}, \tilde{\mathcal{M}}'$  are at least  $\varepsilon_1^{c_1}$ .

Now we can apply Lemma 9.9 on  $\tilde{\mathcal{M}}, \tilde{\mathcal{M}}'$  and get partitions of components of  $\tilde{\mathcal{M}}, \tilde{\mathcal{M}}'$ . For  $i \in [\ell]$ , let  $\mathcal{M}_i$  and  $\mathcal{M}'_i$  be the mixtures defined in Lemma 9.9. We can apply a linear transformation to make  $\frac{1}{2}\mathcal{M}_i + \frac{1}{2}\mathcal{M}'_i$  in isotropic position. Since the total variation distance is invariant under linear transformations, so we still have both conclusions in Lemma 9.9. Let  $\varepsilon_2 = \varepsilon_1^{c_2}$ . Then  $d_{\text{TV}}(\mathcal{M}_i, \mathcal{M}_{\pi(i)}) \leq O(\varepsilon_2)$  and  $\frac{1}{2}\mathcal{M} + \frac{1}{2}\mathcal{M}'$  satisfies Lemma 9.4 with  $\delta = \varepsilon_2^{c_3}$ . Weights of both mixtures increase when we do the partition. So minimum weights are at least  $\varepsilon_1^{c_1} \geq \varepsilon_1^{c_2 c_3} = \varepsilon_2^{c_3}$ .

We now prove the statement on these smaller mixtures. First we can use Lemma 9.10 to merge close parameters of  $\mathcal{M}_i, \mathcal{M}'_i$  so that all pairs of parameters are either equal or separated by  $\varepsilon_2^{c_4 c_5}$ . Under this simplification, Lemma 9.11 shows that there is a perfect matching between the same

components in two mixtures and their weights are almost the same. By the last statement in Lemma 9.10, it is also a matching between components of  $\mathcal{M}_i$  and  $\mathcal{M}'_i$  by combining  $\pi$  and  $\pi_1, \pi_2$ . Moreover, if  $\tilde{G}_j = \tilde{G}'_{\pi(j)}$ , then  $d_{\text{TV}}(G_\ell, G'_{\ell'}) \leq \text{poly}(\varepsilon_2)$  for all  $\ell, \ell'$  such that  $\pi_1(\ell) = j, \pi_2(\ell') = \pi(j)$ . Repeating the argument for all pieces in  $\tilde{\mathcal{M}}, \tilde{\mathcal{M}}'$  completes the proof.  $\square$

## 9.1 Proof of Lemma 9.9

In this section, we will prove Lemma 9.9. The following fact in [Liu-Moitra] shows that a good set of clusters of one mixture exists.

**Fact 9.12** (Claim 7.6 in [LM'20]). *Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  be a mixture of Gaussians. For any constants  $0 < \delta < 1$  and  $\varepsilon > 0$ , there exists  $t \in [k^2]$  such that there exists a partition (possibly trivial) of  $[k]$  into sets  $R_1, \dots, R_\ell$  such that*

1. *If we draw edges between all pairs  $i, j$  such that  $d_{\text{TV}}(G_i, G_j) \leq 1 - \varepsilon^{\delta^t}$ , then each piece of the partition is connected*
2. *For any  $i, j$  in different pieces of the partition,  $d_{\text{TV}}(G_i, G_j) \geq 1 - \varepsilon^{\delta^{t-1}}$ .*

*Remark.* Fact 9.12 can be applied to a set of Gaussians instead of a mixture of Gaussians by randomly assigning positive weights for all Gaussians.

**Lemma 9.13.** *For any constant  $0 < c < 1$ , suppose  $\mathcal{M} = \sum_{i=1}^{k_1} w_i A_i, \mathcal{M}' = \sum_{i=1}^{k_2} w'_i B_i$  are two mixtures of arbitrary distributions with  $d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \leq \varepsilon$  and  $w_i, w'_i \geq \varepsilon^c$ . If for any  $i \neq j$ ,  $d_{\text{TV}}(A_i, B_j) \geq 1 - \varepsilon$ , then  $k_1 = k_2$  and  $d_{\text{TV}}(A_i, B_i) \leq \text{poly}_{k_1}(\varepsilon)$  for all  $i \in [k_1]$ .*

*Proof.* Suppose  $\pi$  is any coupling of  $\mathcal{M}$  and  $\mathcal{M}'$  and  $X, Y$  are random variables with distributions  $\mathcal{M}$  and  $\mathcal{M}'$ . Then  $d_{\text{TV}}(\mathcal{M}, \mathcal{M}') = \min_{\pi} \{\mathbb{P}_{\pi}(X \neq Y)\}$ . We define  $\pi$  to be the optimal coupling such that  $d_{\text{TV}}(\mathcal{M}, \mathcal{M}') = \mathbb{P}_{\pi}(X \neq Y)$ . Then we can define  $\hat{\pi}$  on variables  $i, j, X, Y$  such that  $\sum_{i \in [k_1], j \in [k_2]} \hat{\pi}(i, j, X, Y) = \pi(X, Y)$  and the marginal distribution  $\hat{\pi}_X$  with fixed  $i$  of  $X$  is  $w_i A_i$  for all  $i \in [k_1]$  and the marginal distribution  $\hat{\pi}_Y$  with fixed  $j$  is  $w'_j B_j$  for all  $j \in [k_2]$ . Let  $P_{ij} = \int_{X, Y} \hat{\pi}(i, j, X, Y) dX dY$  and  $A_{ij} = \frac{1}{P_{ij}} \int_Y \hat{\pi}(i, j, X, Y) dY$  be distributions on  $X$ ,  $B_{ij} = \frac{1}{P_{ij}} \int_X \hat{\pi}(i, j, X, Y) dX$  be distributions on  $Y$ . Then we have

$$\begin{aligned} d_{\text{TV}}(\mathcal{M}, \mathcal{M}') &= \mathbb{P}_{\pi}(X \neq Y) = \mathbb{P}_{\hat{\pi}}(X \neq Y) \\ &= \sum_{i, j} P_{ij} \mathbb{P}_{\hat{\pi}}(X \neq Y \mid i, j) \\ &\geq \sum_{i, j} P_{ij} \cdot d_{\text{TV}}(A_{ij}, B_{ij}). \end{aligned} \tag{9.2}$$

By the definition of  $P_{ij}, A_{ij}, B_{ij}$ ,

$$w_i A_i = P_{ij} A_{ij} + \sum_{j' \neq j} P_{ij'} A_{ij'}$$

$$w'_i B_i = P_{ij} B_{ij} + \sum_{i' \neq i} P_{i'j} B_{i'j}$$

Dividing both sides by  $\max\{w_i, w'_j\}$ , we get

$$A_i = \frac{P_{ij}}{\max\{w_i, w'_j\}} A_{ij} + \left(1 - \frac{w_i}{\max\{w_i, w'_j\}}\right) A_i + \sum_{j' \neq j} \frac{P_{ij'}}{\max\{w_i, w'_j\}} A_{ij'}$$

$$B_i = \frac{P_{ij}}{\max\{w_i, w'_j\}} B_{ij} + \left(1 - \frac{w'_j}{\max\{w_i, w'_j\}}\right) B_i + \sum_{i' \neq i} \frac{P_{i'j}}{\max\{w_i, w'_j\}} B_{i'j}$$

From the above two equations, we can write  $A_i, B_i$  as linear combinations of two distributions.

$$A_i = \frac{P_{ij}}{\max\{w_i, w'_j\}} A_{ij} + \left(1 - \frac{P_{ij}}{\max\{w_i, w'_j\}}\right) A'_i$$

$$B_i = \frac{P_{ij}}{\max\{w_i, w'_j\}} B_{ij} + \left(1 - \frac{P_{ij}}{\max\{w_i, w'_j\}}\right) B'_i$$

Then by the triangle inequality,

$$d_{\text{TV}}(A_i, B_i) \leq \frac{P_{ij}}{\max\{w_i, w'_j\}} d_{\text{TV}}(A_{ij}, B_{ij}) + \left(1 - \frac{P_{ij}}{\max\{w_i, w'_j\}}\right) d_{\text{TV}}(A'_i, B'_i)$$

$$P_{ij} \cdot d_{\text{TV}}(A_{ij}, B_{ij}) \geq P_{ij} - (1 - d_{\text{TV}}(A_i, B_i)) \max\{w_i, w'_j\}. \quad (9.3)$$

Combining (9.2) and (9.3), we have the following inequality on the TV distance between mixtures and the TV distance between components

$$d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \geq \sum_{i,j} \left( P_{ij} - (1 - d_{\text{TV}}(A_i, B_j)) \max\{w_i, w'_j\} \right). \quad (9.4)$$

By the lower bounds on  $d_{\text{TV}}(A_i, B_j)$ , we have

$$\begin{aligned} \varepsilon &\geq d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \geq \sum_{i,j} P_{ij} - \sum_{i \neq j} (1 - d_{\text{TV}}(A_i, B_j)) \max\{w_i, w'_j\} - \sum_i (1 - d_{\text{TV}}(A_i, B_i)) \max\{w_i, w'_i\} \\ &\geq 1 - \sum_{i \neq j} \varepsilon - \sum_i \max\{w_i, w'_i\} + \sum_i \max\{w_i, w'_i\} d_{\text{TV}}(A_i, B_i) \\ &\geq 1 - \sum_{i \neq j} \varepsilon - \sum_i \max\{w_i, w'_i\} + w_{\min} d_{\text{TV}}(A_1, B_1) \end{aligned} \quad (9.5)$$

where  $A_1, B_1$  can be replaced by any  $A_i, B_i$  pair. Let  $k = \max\{k_1, k_2\}$ . When  $i \neq j$  and  $d_{\text{TV}}(A_i, B_j) \geq 1 - \varepsilon$ , we plug it into Equation (9.4) and get

$$\varepsilon \geq d_{\text{TV}}(\mathcal{M}, \mathcal{M}') \geq \sum_{i \neq j} P_{ij} - (1 - d_{\text{TV}}(A_i, B_j)) \max\{w_i, w'_j\} \geq \sum_{i \neq j} (P_{ij} - \varepsilon).$$

This implies

$$\sum_{i \neq j} P_{ij} \leq k^2 \varepsilon.$$

Then we can bound  $\sum_i \max\{w_i, w'_i\} - 1$  in Equation (9.5)

$$\sum_i \max\{w_i, w'_i\} - 1 = \sum_i \max\{w_i, w'_i\} - w_i \leq \sum_{i \neq j} P_{ij} \leq k^2 \varepsilon.$$

Plugging this bound into Equation (9.5), for any  $i$ , we have

$$d_{\text{TV}}(A_i, B_i) \leq \frac{1}{w_{\min}} \left( k^2 \varepsilon + \sum_i \max\{w_i, w'_i\} - 1 \right) \leq \frac{2k^2 \varepsilon}{\varepsilon^c}.$$

□

*Proof of Lemma 9.9.* We apply Fact 9.12 on the union set of components of  $\mathcal{M}$  and  $\mathcal{M}'$  with parameter  $\delta$  to find a partition  $R_1, \dots, R_\ell$ . Let

$$\begin{aligned} \mathcal{M}_i &= \frac{\sum_{G_j \in R_i} w_j G_j}{\sum_{G_j \in R_i} w_j} \\ \mathcal{M}'_i &= \frac{\sum_{G'_j \in R_i} w'_j G'_j}{\sum_{G'_j \in R_i} w'_j}. \end{aligned}$$

Then for any  $i \neq j$ , we know  $d_{\text{TV}}(G_a, G'_b) \geq 1 - \varepsilon^{\delta^{t-1}}$  for  $G_a \in R_i, G'_b \in R_j$ . By (9.4) in the proof of Lemma 9.13, we have

$$d_{\text{TV}}(\mathcal{M}_i, \mathcal{M}'_j) \geq 1 - 2k \varepsilon^{\delta^{t-1}}.$$

Then by Lemma 9.13, for any  $i$ , there exists  $a$  such that  $d_{\text{TV}}(\mathcal{M}_i, \mathcal{M}'_i) \leq \varepsilon^{a \delta^{t-1}}$ . Let  $c_2 = a \delta^{t-1}$ . If we set  $\delta = c_2 c_3 / \delta^{t-1} = a c_3$ , the partition satisfies the second conclusion. □

## Acknowledgments

We thank an anonymous reviewer for pointing out an issue with a technical statement (Fact 2.35 in the previous arXiv version, replaced by Lemmas 2.40 and 2.43 in the current version) that claimed a bound on the variance of a more general class of distributions than what is needed in our approach.

A.B. was supported by the Office of Naval Research (ONR) grant N00014-18-1-2562, and the National Science Foundation (NSF) Grant No. CCF-1815840. I.D. was supported by NSF Award CCF-1652862 (CAREER), a Sloan Research Fellowship, and a DARPA Learning with Less Labels (LwLL) grant. H.J. and S.S.V. were supported in part by NSF awards AF-1909756 and AF-2007443. D.M.K. was supported by NSF Award CCF-1553288 (CAREER) and a Sloan Research Fellowship. Part of this work was done while A.B. was visiting the Simons Institute for the Theory of Computing.

## References

- [AK01] S. Arora and R. Kannan, *Learning mixtures of arbitrary Gaussians*, Proceedings of the 33rd Symposium on Theory of Computing, 2001, pp. 247–257. [1](#)
- [AM05] D. Achlioptas and F. McSherry, *On spectral learning of mixtures of distributions*, Proceedings of the Eighteenth Annual Conference on Learning Theory (COLT), 2005, pp. 458–469. [1](#)
- [BBV08] M.-F. Balcan, A. Blum, and S. Vempala, *A discriminative framework for clustering via similarity functions*, Proceedings of the fortieth annual ACM symposium on Theory of computing, 2008, pp. 671–680. [7](#)
- [BDJ<sup>+</sup>20] A. Bakshi, I. Diakonikolas, H. Jia, D. M. Kane, P. K. Kothari, and S. S. Vempala, *Robustly learning mixtures of  $k$  arbitrary gaussians*, CoRR [abs/2012.02119](#) (2020), Available at <https://arxiv.org/abs/2012.02119v2>. [3](#), [6](#), [7](#)
- [BDLS17] S. Balakrishnan, S. S. Du, J. Li, and A. Singh, *Computationally efficient robust sparse estimation in high dimensions*, Proc. 30th Annual Conference on Learning Theory, 2017, pp. 169–212. [7](#)
- [BK20a] A. Bakshi and P. Kothari, *List-decodable subspace recovery via sum-of-squares*, arXiv preprint arXiv:2002.05139 (2020). [7](#), [13](#), [22](#), [26](#)
- [BK20b] A. Bakshi and P. Kothari, *Outlier-robust clustering of non-spherical mixtures*, CoRR [abs/2005.02970](#) (2020). [2](#), [4](#), [5](#), [7](#), [9](#), [10](#), [21](#), [24](#), [26](#), [44](#), [45](#), [49](#), [50](#), [51](#), [71](#)
- [BKS15] B. Barak, J. A. Kelner, and D. Steurer, *Dictionary learning and tensor decomposition via the sum-of-squares method [extended abstract]*, STOC’15—Proceedings of the 2015 ACM Symposium on Theory of Computing, ACM, New York, 2015, pp. 143–151. MR 3388192 [7](#), [12](#), [21](#)
- [BP20] A. Bakshi and A. Prasad, *Robust linear regression: Optimal rates in polynomial time*, arXiv preprint arXiv:2007.01394 (2020). [7](#)
- [BS10] M. Belkin and K. Sinha, *Polynomial learning of distribution families*, FOCS, 2010, pp. 103–112. [1](#), [4](#)
- [BS16] B. Barak and D. Steurer, *Proofs, beliefs, and algorithms through the lens of sum-of-squares*, 2016, Lecture notes in preparation, available on <http://sumofsquares.org>. [18](#)
- [BV08] S. C. Brubaker and S. Vempala, *Isotropic PCA and Affine-Invariant Clustering*, Proc. 49th IEEE Symposium on Foundations of Computer Science, 2008, pp. 551–560. [1](#)
- [CAT<sup>+</sup>20] Y. Cherapanamjeri, E. Aras, N. Tripuraneni, M. I. Jordan, N. Flammarion, and P. L. Bartlett, *Optimal robust linear regression in nearly linear time*, arXiv preprint arXiv:2007.08137 (2020). [7](#)

- [CDGS20] Y. Cheng, I. Diakonikolas, R. Ge, and M. Soltanolkotabi, *High-dimensional robust mean estimation via gradient descent*, CoRR **abs/2005.01378** (2020). [7](#)
- [CDKS18] Y. Cheng, I. Diakonikolas, D. M. Kane, and A. Stewart, *Robust learning of fixed-structure Bayesian networks*, Proc. 33rd Annual Conference on Neural Information Processing Systems (NeurIPS), 2018, pp. 10304–10316. [7](#)
- [CSV17] M. Charikar, J. Steinhardt, and G. Valiant, *Learning from untrusted data*, Proc. 49th Annual ACM Symposium on Theory of Computing, 2017, pp. 47–60. [7](#)
- [CW01] A. Carbery and J. Wright, *Distributional and  $L^q$  norm inequalities for polynomials over convex bodies in  $R^n$* , Mathematical Research Letters **8** (2001), no. 3, 233–248. [18](#)
- [Das99] S. Dasgupta, *Learning mixtures of Gaussians*, Proceedings of the 40th Annual Symposium on Foundations of Computer Science, 1999, pp. 634–644. [1](#)
- [DDS12] C. Daskalakis, I. Diakonikolas, and R.A. Servedio, *Learning Poisson Binomial Distributions*, Proceedings of the 44th Symposium on Theory of Computing, 2012, pp. 709–728. [28](#)
- [DDS15] A. De, I. Diakonikolas, and R. Servedio, *Learning from satisfying assignments*, Proceedings of the 26th Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015, 2015, pp. 478–497. [28](#)
- [DGJ<sup>+</sup>10] I. Diakonikolas, P. Gopalan, R. Jaiswal, R. Servedio, and E. Viola, *Bounded independence fools halfspaces*, SIAM J. on Comput. **39** (2010), no. 8, 3441–3462. [22](#)
- [DHKK20] I. Diakonikolas, S. B. Hopkins, D. Kane, and S. Karmalkar, *Robustly learning any clusterable mixture of gaussians*, CoRR **abs/2005.06417** (2020). [2](#), [4](#), [5](#), [7](#), [9](#), [10](#), [44](#)
- [DK14] C. Daskalakis and G. Kamath, *Faster and sample near-optimal algorithms for proper learning mixtures of Gaussians*, Proc. 27th Annual Conference on Learning Theory (COLT), 2014, pp. 1183–1213. [28](#)
- [DK19] I. Diakonikolas and D. M. Kane, *Recent advances in algorithmic high-dimensional robust statistics*, CoRR **abs/1911.05911** (2019). [1](#), [7](#)
- [DKK<sup>+</sup>16] I. Diakonikolas, G. Kamath, D. M. Kane, J. Li, A. Moitra, and A. Stewart, *Robust estimators in high dimensions without the computational intractability*, Proc. 57th IEEE Symposium on Foundations of Computer Science (FOCS), 2016, pp. 655–664. [1](#), [2](#), [7](#), [26](#), [28](#)
- [DKK<sup>+</sup>17] I. Diakonikolas, G. Kamath, D. M. Kane, J. Li, A. Moitra, and A. Stewart, *Being robust (in high dimensions) can be practical*, Proc. 34th International Conference on Machine Learning (ICML), 2017, pp. 999–1008. [7](#), [23](#)
- [DKK<sup>+</sup>18] I. Diakonikolas, G. Kamath, D. M. Kane, J. Li, J. Steinhardt, and A. Stewart, *Sever: A robust meta-algorithm for stochastic optimization*, CoRR **abs/1803.02815** (2018), Conference version in ICML 2019. [7](#)



- [DKK<sup>+</sup>19] I. Diakonikolas, S. Karmalkar, D. Kane, E. Price, and A. Stewart, *Outlier-robust high-dimensional sparse estimation via iterative filtering*, Advances in Neural Information Processing Systems 33, NeurIPS 2019, 2019. [7](#)
- [DKS17] I. Diakonikolas, D. M. Kane, and A. Stewart, *Statistical query lower bounds for robust estimation of high-dimensional Gaussians and Gaussian mixtures*, Proc. 58th IEEE Symposium on Foundations of Computer Science (FOCS), 2017, pp. 73–84. [3](#), [8](#), [11](#)
- [DKS18] I. Diakonikolas, D. M. Kane, and A. Stewart, *List-decodable robust mean estimation and learning mixtures of spherical Gaussians*, Proc. 50th Annual ACM Symposium on Theory of Computing (STOC), 2018, pp. 1047–1060. [2](#), [111](#)
- [DKS19] I. Diakonikolas, W. Kong, and A. Stewart, *Efficient algorithms and lower bounds for robust linear regression*, Proc. 30th Annual Symposium on Discrete Algorithms (SODA), 2019, pp. 2745–2754. [7](#)
- [DL01] L. Devroye and G. Lugosi, *Combinatorial methods in density estimation*, Springer Series in Statistics, Springer, 2001. [28](#), [112](#)
- [DRST14] I. Diakonikolas, P. Raghavendra, R. A. Servedio, and L. Y. Tan, *Average sensitivity and noise sensitivity of polynomial threshold functions*, SIAM J. Comput. **43** (2014), no. 1, 231–253. [104](#)
- [DVW19] I. Diakonikolas, S. Vempala, and D. Woodruff, *Research vignette: Foundations of data science*, UC Berkeley Simons Institute newsletter (2019). [2](#)
- [EL73] Paul Erdős and László Lovász, *Problems and results on 3-chromatic hypergraphs and some related questions*, Colloquia Mathematica Societatis Janos Bolyai 10. Infinite and Finite Sets, Keszthely (Hungary), Citeseer, 1973. [115](#)
- [FJK96] A. Frieze, M. Jerrum, and R. Kannan, *Learning linear transformations*, focs1996, 1996, pp. 359–368. [8](#)
- [FKP19] N. Fleming, P. Kothari, and T. Pitassi, *Semialgebraic proofs and efficient algorithm design*, Foundations and Trends® in Theoretical Computer Science **14** (2019), no. 1-2, 1–221. [7](#), [18](#), [49](#), [71](#)
- [GHK15] R. Ge, Q. Huang, and S. M. Kakade, *Learning mixtures of gaussians in high dimensions*, Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, 2015, pp. 761–770. [8](#)
- [GKS20] A. Garg, N. Kayal, and C. Saha, *Learning sums of powers of low-degree polynomials in the non-degenerate case*, CoRR [abs/2004.06898](#) (2020). [11](#)
- [GLS81] M. Grötschel, L. Lovász, and A. Schrijver, *The ellipsoid method and its consequences in combinatorial optimization*, Combinatorica **1** (1981), no. 2, 169–197. MR 625550 [19](#)

- [GVX14] N. Goyal, S. Vempala, and Y. Xiao, *Fourier pca and robust tensor decomposition*, Proceedings of the forty-sixth annual ACM symposium on Theory of computing, 2014, pp. 584–593. [8](#)
- [Her20] Hermite Polynomials, *Inverse explicit expressions* — *Wikipedia, the free encyclopedia*, 2020, [Online]. [17](#)
- [HK13] D. Hsu and S. M. Kakade, *Learning mixtures of spherical gaussians: moment methods and spectral decompositions*, Innovations in Theoretical Computer Science, ITCS '13, 2013, pp. 11–20. [8](#)
- [HL18] S. B. Hopkins and J. Li, *Mixture models, robustness, and sum of squares proofs*, Proc. 50th Annual ACM Symposium on Theory of Computing (STOC), 2018, pp. 1021–1034. [2](#), [7](#), [71](#)
- [HP15] M. Hardt and E. Price, *Tight bounds for learning a mixture of two gaussians*, Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, 2015, pp. 753–760. [1](#)
- [HR09] P. J. Huber and E. M. Ronchetti, *Robust statistics*, Wiley New York, 2009. [1](#)
- [HRRS86] F. R. Hampel, E. M. Ronchetti, P. J. Rousseeuw, and W. A. Stahel, *Robust statistics. the approach based on influence functions*, Wiley New York, 1986. [1](#)
- [Hub64] P. J. Huber, *Robust estimation of a location parameter*, Ann. Math. Statist. **35** (1964), no. 1, 73–101. [2](#)
- [Kan20] D. M. Kane, *Robust learning of mixtures of gaussians*, CoRR **abs/2007.05912** (2020). [2](#), [4](#), [5](#), [10](#), [11](#), [16](#), [17](#), [28](#)
- [KKK19] S. Karmalkar, A. Klivans, and P. Kothari, *List-decodable linear regression*, Advances in Neural Information Processing Systems, 2019, pp. 7423–7432. [7](#), [13](#), [22](#), [26](#)
- [KKM18] A. Klivans, P. Kothari, and R. Meka, *Efficient algorithms for outlier-robust regression*, Proc. 31st Annual Conference on Learning Theory (COLT), 2018, pp. 1420–1430. [7](#)
- [KMV10] A. T. Kalai, A. Moitra, and G. Valiant, *Efficiently learning mixtures of two Gaussians*, STOC, 2010, pp. 553–562. [1](#), [9](#)
- [KOTZ14] M. Kauters, R. O’Donnell, L. Tan, and Y. Zhou, *Hypercontractive inequalities via sos, and the frankl-rödl graph*, SODA, SIAM, 2014, pp. 1644–1658. [23](#), [108](#)
- [KS17a] P. K. Kothari and J. Steinhardt, *Better agnostic clustering via relaxed tensor norms*, CoRR **abs/1711.07465** (2017). [7](#), [71](#)
- [KS17b] P. K. Kothari and D. Steurer, *Outlier-robust moment-estimation via sum-of-squares*, CoRR **abs/1711.11581** (2017). [7](#), [10](#), [22](#), [23](#), [24](#), [26](#), [71](#)
- [KSS18] P. K. Kothari, J. Steinhardt, and D. Steurer, *Robust moment estimation and improved clustering via sum of squares*, Proc. 50th Annual ACM Symposium on Theory of Computing (STOC), 2018, pp. 1035–1046. [2](#)

- [KSV08] R. Kannan, H. Salmasian, and S. Vempala, *The spectral method for general mixture models*, SIAM J. Comput. **38** (2008), no. 3, 1141–1156. [1](#)
- [Las01] J. B. Lasserre, *New positive semidefinite relaxations for nonconvex quadratic programs*, Advances in convex analysis and global optimization (Pythagorion, 2000), Nonconvex Optim. Appl., vol. 54, Kluwer Acad. Publ., Dordrecht, 2001, pp. 319–331. MR 1846160 [19](#)
- [LM20] A. Liu and A. Moitra, *Settling the robust learnability of mixtures of gaussians*, arXiv preprint arXiv:2011.03622 (2020). [3](#), [4](#), [5](#), [6](#), [7](#), [15](#), [89](#), [90](#), [93](#)
- [LRV16] K. A. Lai, A. B. Rao, and S. Vempala, *Agnostic estimation of mean and covariance*, Proc. 57th IEEE Symposium on Foundations of Computer Science (FOCS), 2016, pp. 665–674. [1](#), [2](#), [7](#)
- [MV10] A. Moitra and G. Valiant, *Settling the polynomial learnability of mixtures of Gaussians*, FOCS, 2010, pp. 93–102. [1](#), [4](#), [9](#), [11](#), [15](#), [39](#), [42](#), [43](#), [81](#)
- [Nes00] Y. Nesterov, *Squared functional systems and optimization problems*, High performance optimization, Appl. Optim., vol. 33, Kluwer Acad. Publ., Dordrecht, 2000, pp. 405–440. MR 1748764 [19](#)
- [O’D14] R. O’Donnell, *Analysis of Boolean functions*, Cambridge University Press, New York, 2014. MR 3443800 [23](#)
- [Par13] P. A. Parrilo, *Polynomial optimization, sums of squares, and applications*, Semidefinite optimization and convex algebraic geometry, MOS-SIAM Ser. Optim., vol. 13, SIAM, Philadelphia, PA, 2013, pp. 47–157. MR 3050242 [19](#)
- [Pea94] K. Pearson, *Contribution to the mathematical theory of evolution*, Phil. Trans. Roy. Soc. A **185** (1894), 71–110. [1](#)
- [PSBR18] A. Prasad, A. S. Suggala, S. Balakrishnan, and P. Ravikumar, *Robust estimation via robust gradient estimation*, arXiv preprint arXiv:1802.06485 (2018). [7](#)
- [RY20a] P. Raghavendra and M. Yau, *List decodable learning via sum of squares*, Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SIAM, 2020, pp. 161–180. [7](#), [13](#), [22](#)
- [RY20b] P. Raghavendra and M. Yau, *List decodable subspace recovery*, arXiv preprint arXiv:2002.03004 (2020). [7](#), [13](#)
- [Sho87] N. Z. Shor, *Quadratic optimization problems*, Izv. Akad. Nauk SSSR Tekhn. Kibernet. (1987), no. 1, 128–139, 222. MR 939596 [19](#)
- [TLM18] B. Tran, J. Li, and A. Madry, *Spectral signatures in backdoor attacks*, Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 2018, pp. 8011–8021. [7](#)

- [VW02] S. Vempala and G. Wang, *A spectral algorithm for learning mixtures of distributions*, Proc. 43rd IEEE Symposium on Foundations of Computer Science (FOCS), 2002, pp. 113–122. [1](#)
- [Yat85] Y. G. Yatracos, *Rates of convergence of minimum distance estimators and Kolmogorov's entropy*, Annals of Statistics **13** (1985), 768–774. [28](#)
- [ZJS20] B. Zhu, J. Jiao, and J. Steinhardt, *Robust estimation via generalized quasi-gradients*, arXiv preprint arXiv:2005.14073 (2020). [7](#)

## A Omitted Proofs

In this subsection, we provide the proofs that were omitted from Section 2 and Section 6.

### A.1 Omitted Proofs from Section 2.1

**Lemma A.1** (Concentration of low-degree polynomials, Lemma 2.9 restated). *Let  $T$  be a  $d$ -dimensional, degree-4 tensor such that  $\|T\|_F \leq \Delta$  for some  $\Delta > 0$  and let  $x, y \sim \mathcal{N}(0, I)$ . Then, with probability at least  $1 - 1/\text{poly}(d)$ , the following holds:*

$$\|T(\cdot, \cdot, x, y)\|_F^2 \leq O(\log(d)\Delta^2) .$$

*Proof.* We note that

$$\begin{aligned} \mathbb{E} \left[ \|T(\cdot, \cdot, x, y)\|_F^2 \right] &= \mathbb{E} \left[ \sum_{i_1, i_2} \left( \sum_{i_3, i_4} T(i_1, i_2, i_3, i_4) x(i_3) y(i_4) \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{i_1, i_2} \left( \sum_{i_3, i_4} T(i_1, i_2, i_3, i_4)^2 x(i_3)^2 y(i_4)^2 \right) \right] \\ &= \sum_{i_1, i_2, i_3, i_4} T(i_1, i_2, i_3, i_4)^2 \leq \Delta^2 . \end{aligned}$$

The second equality follows from the fact that  $x(i_3), y(i_4)$  are independent and have zero means. So the only non-zero terms are the squares. The third equality follows from the fact that  $x(i_3), y(i_4)$  are independent with unit variances. Observe that  $\|T(\cdot, \cdot, x, y)\|_F^2$  is a degree-2 polynomial in Gaussian random variables. Using standard concentration bounds for low-degree Gaussian polynomials (see, e.g., Theorem 2.3 in [DRST14]), we obtain

$$\mathbb{P} \left[ \|T(\cdot, \cdot, x, y)\|_F^2 \geq t^2 \mathbb{E} \left[ \|T(\cdot, \cdot, x, y)\|_F^2 \right] \right] \leq \exp(-ct) .$$

Setting  $t = \Omega(\log(d))$  completes the proof. □

## A.2 Omitted Proofs from Section 2.2

**Lemma A.2** (Spectral SoS Proofs, Lemma 2.23 restated). *Let  $A$  be a  $d \times d$  matrix. Then for  $d$ -dimensional vector-valued indeterminate  $v$ , we have:*

$$\frac{|v}{2} \{v^\top A v \leq \|A\|_2 \|v\|_2^2\}.$$

*Proof.* Note that  $v$  is the only variable in the proof here ( $A$  is a matrix of constants). We note that  $A \leq \|A\|_2 I$  or  $\|A\|_2 I - A$  is PSD and thus  $\|A\|_2 I - A = QQ^\top$  for some  $d \times d$  matrix  $Q$ . Thus,  $\|Qv\|_2^2 = v^\top (\|A\|_2 I - A)v = \|A\|_2 \|v\|_2^2 - v^\top A v$ . Thus,  $\|A\|_2 \|v\|_2^2 - v^\top A v$  is a sum of squares polynomial (namely  $\|Qv\|_2^2$ ) in variable  $v$ . This completes the proof.  $\square$

**Lemma A.3** (Frobenius Norms of Products of Matrices, Lemma 2.25 restated). *Let  $B$  be a  $d \times d$  matrix valued indeterminate for some  $d \in \mathbb{N}$ . Then, for any  $0 \leq A \leq I$ ,*

$$\frac{|B}{2} \{\|AB\|_F^2 \leq \|B\|_F^2\},$$

and,

$$\frac{|B}{2} \{\|BA\|_F^2 \leq \|B\|_F^2\},$$

*Proof.* The proof of the second claim is similar so we prove only the first. We have:

$$\frac{|B}{2} \{\|B\|_F^2 = \|(A + I - A)B\|_F^2 = \|AB\|_F^2 + \|(I - A)B\|_F^2 + 2\text{tr}((I - A)BB^\top A)\}$$

Now,  $A - A^2 \geq 0$ , thus,  $A - A^2 = RR^\top$  for some  $d \times d$  matrix  $R$ . Thus,  $\text{tr}((A - A^2)BB^\top) = \text{tr}(RR^\top BB^\top) = \|BR\|_F^2$  - a sum of squares polynomial of degree 2 in indeterminate  $B$ . Thus,  $\frac{|B}{2} \{\text{tr}((A - A^2)BB^\top) \geq 0\}$ .  $\square$

## A.3 Omitted Proofs from Section 2.3

**Lemma A.4** (Shifts Cannot Decrease Variance, Lemma 2.32 restated). *Let  $\mathcal{D}$  be a distribution on  $\mathbb{R}^d$ ,  $Q$  be a  $d \times d$  matrix-valued indeterminate, and  $C$  be a scalar-valued indeterminate. Then, we have that*

$$\frac{|Q,C}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x) - \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] \right)^2 \right] \leq \mathbb{E}_{x \sim \mathcal{D}} \left[ (Q(x) - C)^2 \right] \right\}.$$

*Proof.*

$$\begin{aligned} \frac{|Q,C}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ (Q(x) - C)^2 \right] \right\} &= \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x) - \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] + \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] - C \right)^2 \right] \\ &= \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x) - \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] \right)^2 \right] + \mathbb{E}_{x \sim \mathcal{D}} \left[ (Q(x) - C)^2 \right] \\ &\quad + 2 \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x) - \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] \right) \left( \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] - C \right) \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x) - \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] \right)^2 \right] + \mathbb{E}_{x \sim \mathcal{D}} [(Q(x) - C)^2] \\
&\geq \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x) - \mathbb{E}_{x \sim \mathcal{D}} [Q(x)] \right)^2 \right] \left. \right\}.
\end{aligned}$$

□

**Lemma A.5** (Shifts of Certifiably Hypercontractive Distributions, Lemma 2.33 restated). *Let  $x$  be a mean-0 random variable with distribution  $\mathcal{D}$  on  $\mathbb{R}^d$  with  $t$ -certifiably  $C$ -hypercontractive degree-2 polynomials. Then, for any fixed constant vector  $c \in \mathbb{R}^d$ , the random variable  $x + c$  also has  $t$ -certifiable  $4C$ -hypercontractive degree-2 polynomials.*

*Proof.* Observe that using that  $\mathbb{E}_{x \sim \mathcal{D}} [x] = 0$ , we have that

$$\frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} [(x+c)^\top Q(x+c)] = \mathbb{E}_{x \sim \mathcal{D}} [x^\top Qx + c^\top Qc] \right\}.$$

Next, by two applications of the SoS Triangle Inequality (Fact 2.21), an application of Lemma 2.32 followed by certifiable hypercontractivity of  $\mathcal{D}$ , we have:

$$\begin{aligned}
&\frac{Q}{t'} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( (x+c)^\top Q(x+c) - \mathbb{E}_{x \sim \mathcal{D}} [(x+c)^\top Q(x+c)] \right)^{t'} \right] \right. \\
&= \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}} [x^\top Qx] \right) + x^\top Qc + c^\top Qx \right)^{t'} \right] \\
&\leq 4^{t'} \left( \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( x^\top Qx - \mathbb{E}_{\mathcal{D}} [x^\top Qx] \right)^{t'} \right] + \mathbb{E}_{x \sim \mathcal{D}} [(x^\top Qc)^{t'}] + \mathbb{E}_{x \sim \mathcal{D}} [(c^\top Qx)^{t'}] \right) \\
&\leq 4^{t'} (Ct')^{t'} \left( \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( x^\top Qx - \mathbb{E}_{\mathcal{D}} [x^\top Qx] \right)^2 \right]^{t'/2} + \mathbb{E}_{x \sim \mathcal{D}} [(x^\top Qc)^2]^{t'/2} + \mathbb{E}_{x \sim \mathcal{D}} [(c^\top Qx)^2]^{t'/2} \right) \left. \right\}.
\end{aligned}$$

On the other hand, notice that

$$\begin{aligned}
&\frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( (x+c)^\top Q(x+c) - \mathbb{E}_{x \sim \mathcal{D}} [(x+c)^\top Q(x+c)] \right)^2 \right] \right. \\
&= \left( \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( x^\top Qx - \mathbb{E}_{\mathcal{D}} [x^\top Qx] \right)^2 \right] + \mathbb{E}_{x \sim \mathcal{D}} [(x^\top Qc)^2] + \mathbb{E}_{x \sim \mathcal{D}} [(c^\top Qx)^2] \right) \left. \right\}.
\end{aligned}$$

Thus,

$$\frac{Q}{t'} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}} [x^\top Qx] \right)^2 \right]^{t'/2} + \left( \mathbb{E}_{x \sim \mathcal{D}} [(x^\top Qc)^2] \right)^{t'/2} + \left( \mathbb{E}_{x \sim \mathcal{D}} [(c^\top Qx)^2] \right)^{t'/2} \right.$$

$$\leq 4^{t'} (Ct')^{t'} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( (x+c)^\top Q(x+c) - \mathbb{E}_{x \sim \mathcal{D}} [(x+c)^\top Q(x+c)] \right)^2 \right]^{t'/2} \right\}.$$

As a result, we obtain:

$$\begin{aligned} & \frac{Q}{t'} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( (x+c)^\top Q(x+c) - \mathbb{E}_{x \sim \mathcal{D}} [(x+c)^\top Q(x+c)] \right)^{t'} \right] \right\} \\ & \leq (4Ct')^{t'} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( (x+c)^\top Q(x+c) - \mathbb{E}_{x \sim \mathcal{D}} [(x+c)^\top Q(x+c)] \right)^2 \right]^{t'/2} \right\}, \end{aligned}$$

which completes the proof.  $\square$

**Lemma A.6** (Mixtures of Certifiably Hypercontractive Distributions, Lemma 2.34 restated). *Let  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$  have  $t$ -certifiable  $C$ -hypercontractive degree-2 polynomials on  $\mathbb{R}^d$ , for some fixed constant  $C$ . Then, any mixture  $\mathcal{D} = \sum_i w_i \mathcal{D}_i$  also has  $t$ -certifiably  $(C/\alpha)$ -hypercontractive degree-2 polynomials for  $\alpha = \min_{i \leq k, w_i > 0} w_i$ .*

*Proof.* Applying Lemma 2.21 followed by SoS Hölder's inequality on the second term and followed by a final application of SoS Hölder's inequality (Fact 2.20), we obtain:

$$\begin{aligned} \frac{Q}{t'} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}} [x^\top Qx] \right)^{t'} \right] \right\} &= \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( x^\top Qx - \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Qx] \right)^{t'} \right] \\ &= \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \left( x^\top Qx - \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Qx] \right)^{t'} \right] \\ &\leq 2^{t'} \left( \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Qx] \right)^{t'} \right] \right) \\ &\quad + \left( \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left[ x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Qx] \right]^{t'} \right) \\ &\leq 2^{t'} \left( (Ct')^{t'} \left( \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Qx] \right)^2 \right]^{t'/2} \right) \right) \\ &\quad + \sum_i w_i \left( \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Qx] \right)^2 \right]^{t'} \right) \\ &\leq \left( \frac{4Ct'}{\alpha} \right)^{t'} \left\{ \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Qx] \right)^2 \right]^{t'/2} \right\}. \end{aligned}$$

On the other hand, note that by Lemma 2.32, we know that

$$\frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}} [x^\top Qx] \right)^2 \right] \right\} = \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}} [x^\top Qx] \right)^2 \right]$$



$$\geq \sum_i w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \left( x^\top Q x - \mathbb{E}_{x \sim \mathcal{D}_i} [x^\top Q x] \right)^2 \right].$$

Combining the two equations above completes the proof.  $\square$

**Corollary A.7** (Certifiable Hypercontractivity of  $k$ -Mixtures of Gaussians, Corollary 2.35 restated). *Let  $\mathcal{D}$  be a  $k$ -mixture of Gaussians  $\sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  with weights  $w_i \geq \alpha$  for every  $i \in [k]$ . Then,  $\mathcal{D}$  has  $t$ -certifiably  $4/\alpha$ -hypercontractive degree-2 polynomials.*

*Proof.* From [KOTZ14], we know that the standard Gaussian random variable has  $t$ -certifiably 1-hypercontractive degree-2 polynomials. From Fact 2.31, we immediately obtain that for any PSD matrix  $\Sigma$ , the Gaussian  $\mathcal{N}(0, \Sigma)$  also has  $t$ -certifiably 1-hypercontractive degree-2 polynomials. From Lemma 2.33, we obtain that for any  $\mu$ , the Gaussian  $\mathcal{N}(\mu, \Sigma)$  has  $t$ -certifiably 4-hypercontractive degree-2 polynomials. Finally, applying Lemma 2.34 to  $\mathcal{D}_i = \mathcal{N}(\mu_i, \Sigma_i)$  and mixture weights  $w_1, w_2, \dots, w_k$ , yields that  $\mathcal{D} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  has  $t$ -certifiably  $4/\alpha$ -hypercontractive degree-2 polynomials. This completes the proof.  $\square$

**Lemma A.8** (Linear Transformations of Certifiably Bounded-Variance Distributions, Lemma 2.39 restated). *For  $d \in \mathbb{N}$ , let  $x$  be a random variable with distribution  $\mathcal{D}$  on  $\mathbb{R}^d$  such that for  $d \times d$  matrix-valued indeterminate  $Q$ ,  $\left| \frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}} (x^\top Q x - \mathbb{E}_{\mathcal{D}} x^\top Q x)^2 \leq \left\| \Sigma^{1/2} Q \Sigma^{1/2} \right\|_F^2 \right\} \right.$ . Let  $A$  be an arbitrary  $d \times d$  matrix and let  $x' = Ax$  be the random variable with covariance  $\Sigma' = AA^\top$ . Then, we have that*

$$\left| \frac{Q}{2} \left\{ \mathbb{E}_{x' \sim \mathcal{D}'} (x'^\top Q x' - \mathbb{E}_{\mathcal{D}'} x'^\top Q x')^2 \leq \left\| \Sigma'^{1/2} Q \Sigma'^{1/2} \right\|_F^2 \right\} \right|.$$

*Proof.* The covariance of  $x'$  is  $AA^\top = \Sigma'$ , say. Let  $\Sigma'^{1/2}$  be the PSD square root of  $\Sigma'$ . The proof follows by noting that  $x'^\top Q x' = (Ax)^\top Q (Ax) = x^\top (A^\top Q A) x$  and that  $\|A^\top Q A\|_F^2 = \text{tr}(A^\top Q A A^\top Q A) = \text{tr}(A A^\top Q A A^\top Q) = \text{tr}(\Sigma' Q \Sigma' Q) = \text{tr}(\Sigma'^{1/2} Q \Sigma'^{1/2} \Sigma'^{1/2} Q \Sigma'^{1/2}) = \left\| \Sigma'^{1/2} Q \Sigma'^{1/2} \right\|_F^2$ .  $\square$

**Lemma A.9** (Variance of Degree-2 Polynomials of Standard Gaussians, Lemma 2.40 restated). *We have that*

$$\left| \frac{Q}{2} \left\{ \mathbb{E}_{\mathcal{N}(0, I)} \left( x^\top Q x - \mathbb{E}_{\mathcal{N}(0, I)} x^\top Q x \right)^2 \leq 3 \|Q\|_F^2 \right\} \right|.$$

*Proof.* We will view  $xx^\top$  and  $I \in \mathbb{R}^{d \times d}$  as  $d^2$ -dimensional vectors. Consider the matrix  $\mathbb{E}_{x \sim \mathcal{N}(0, I)} (xx^\top - I)(xx^\top - I)^\top$ . The diagonal of this matrix is  $2I_{d^2}$ . The off-diagonal part has exactly one non-zero entry in any row (which corresponds to entry indexed by  $(i, j)$  and  $(j, i)$  for  $i \neq j$ ), and thus has spectral norm at most 1 by the Gershgorin circle theorem. Thus,  $\mathbb{E}_{x \sim \mathcal{N}(0, I)} (xx^\top - I)(xx^\top - I)^\top \leq 3I_{d^2}$ .

We thus have:

$$\begin{aligned} \left| \frac{Q}{2} \left\{ \mathbb{E}_{\mathcal{N}(0, I)} \left( x^\top Q x - \mathbb{E}_{\mathcal{N}(0, I)} x^\top Q x \right)^2 \right. \right. \\ \left. \left. \leq \left\| \mathbb{E}_{x \sim \mathcal{N}(0, I)} (xx^\top - I)(xx^\top - I)^\top \right\|_2 \|Q\|_F^2 \leq 3 \|I_{d^2}\|_2 \|Q\|_F^2 = 3 \|Q\|_F^2 \right\} \right| \quad (\text{A.1}) \end{aligned}$$

□

**Lemma A.10** (Variance of Degree-2 Polynomials of Mixtures, Lemma 2.43 restated). *Let  $\mathcal{M} = \sum_i w_i \mathcal{D}_i$  be a  $k$ -mixture of distributions  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$  with means  $\mu_i$  and covariances  $\Sigma_i$ . Let  $\mu = \sum_i w_i \mu_i$  be the mean of  $\mathcal{M}$ . Suppose that each of  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$  have certifiably  $C$ -bounded-variance i.e. for  $Q$ : a symmetric  $d \times d$  matrix-valued indeterminate.*

$$\frac{|Q|}{2} \left\{ \mathbb{E}_{x' \sim \mathcal{D}_i} (x'^{\top} Q x' - \mathbb{E}_{\mathcal{D}_i} x'^{\top} Q x')^2 \leq C \left\| \Sigma^{1/2} Q \Sigma^{1/2} \right\|_F^2 \right\}.$$

Further, suppose that for some  $H > 1$ ,  $\|\mu_i - \mu\|_2^2, \|\Sigma_i - I\|_F \leq H$  for every  $1 \leq i \leq k$ . Then, we have that

$$\frac{|Q|}{2} \left\{ \mathbb{E}_{x \sim \mathcal{M}} \left[ \left( x^{\top} Q x - \mathbb{E}_{x \sim \mathcal{M}} [x^{\top} Q x] \right)^2 \right] \leq 100CH^2 \|Q\|_F^2 \right\}.$$

*Proof.* We have the following sequence of (in-)equalities:

$$\frac{|Q|}{2} \left\{ \mathbb{E}_{x \sim \mathcal{M}} \left( x^{\top} Q x - \mathbb{E}_{x \sim \mathcal{M}} x^{\top} Q x \right)^2 = \sum_{i \leq k} w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left( x^{\top} Q x - \mathbb{E}_{x \sim \mathcal{M}} x^{\top} Q x \right)^2 \right. \quad (\text{A.2})$$

$$= \sum_{i \leq k} w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left( x^{\top} Q x - \mathbb{E}_{x \sim \mathcal{D}_i} x^{\top} Q x + \mathbb{E}_{x \sim \mathcal{D}_i} x^{\top} Q x - \mathbb{E}_{x \sim \mathcal{M}} x^{\top} Q x \right)^2 \quad (\text{A.3})$$

$$\leq 2 \sum_{i \leq k} w_i \mathbb{E}_{x \sim \mathcal{D}_i} \left( x^{\top} Q x - \mathbb{E}_{x \sim \mathcal{D}_i} x^{\top} Q x \right)^2 + 2 \sum_i w_i \left( \mathbb{E}_{x \sim \mathcal{D}_i} x^{\top} Q x - \mathbb{E}_{x \sim \mathcal{M}} x^{\top} Q x \right)^2 \left. \right\}, \quad (\text{A.4})$$

where the third line follows from Fact 2.21 (SoS Almost Triangle Inequality).

Let us first bound the 2nd term in the RHS above. Towards that, let  $\Sigma = \sum_i w_i ((\mu_i - \mu)(\mu_i - \mu)^{\top} + \Sigma_i)$  be the covariance of the mixture  $\mathcal{M}$ . Then, notice that  $\Sigma = \sum_i w_i ((\mu_i - \mu)(\mu_i - \mu)^{\top} + \Sigma_i) = \sum_i w_i \mu_i \mu_i^{\top} + \sum_i w_i \Sigma_i - \mu \mu^{\top}$ . Thus, we can write

$$\begin{aligned} \mu_i \mu_i^{\top} + \Sigma_i - \Sigma - \mu \mu^{\top} &= \sum_{j \neq i} w_j (\mu_i \mu_i^{\top} - \mu_j \mu_j^{\top}) + \sum_{j \neq i} w_j (\Sigma_i - \Sigma_j) \\ &= \sum_{j \neq i} w_j (\mu_j - \mu)(\mu_j - \mu)^{\top} - \sum_{j \neq i} (\mu - \mu_j)(\mu - \mu_j)^{\top} + \sum_{j \neq i} w_j (\Sigma_i - \Sigma_j) \\ &= \sum_{j \neq i} w_j (\mu_j - \mu)(\mu_j - \mu)^{\top} - \sum_{j \neq i} (\mu - \mu_j)(\mu - \mu_j)^{\top} + \sum_{j \neq i} w_j (\Sigma_i - I) - \sum_{j \neq i} w_j (\Sigma_j - I). \end{aligned}$$

Here, in the second to last step, we added and subtracted  $\sum_{j \neq i} w_j \mu \mu^{\top}$  and used that  $\sum_i w_i \mu_i = \mu$ , and in the last step we added and subtracted  $\sum_{j \neq i} w_j I$ .

By application of the triangle inequality for Frobenius norm to the RHS of the above, we have that:

$$\|\mu_i \mu_i^{\top} + \Sigma_i - \mu \mu^{\top} - \Sigma\|_F \leq \sum_{j \neq i} w_j \|(\mu_i - \mu)(\mu_i - \mu)^{\top}\|_F + \sum_{j \neq i} w_j \|(\mu_j - \mu)(\mu_j - \mu)^{\top}\|_F$$

$$+ \sum_{j \neq i} w_j \|\Sigma_i - I\|_F + \sum_{j \neq i} w_j \|(I - \Sigma_j)\|_F \leq H + H + H + H = 4H.$$

Using the SoS version of the Cauchy-Schwarz inequality (Fact 2.20) on indeterminate  $Q$  and constant  $\mu\mu^\top - \mu_i\mu_i^\top + \Sigma_i - \Sigma$  and the above bound, we have:

$$\begin{aligned} & \left| \frac{Q}{2} \left\{ \sum_i w_i \left( \mathbb{E}_{x \sim \mathcal{D}_i} x^\top Q x - \mathbb{E}_{x \sim \mathcal{M}} x^\top Q x \right)^2 \right. \right. \\ & \quad \left. \left. \leq \sum_i w_i \|\mu\mu^\top - \mu_i\mu_i^\top + \Sigma_i - \Sigma\|_F^2 \|Q\|_F^2 \leq 16H^2 \sum_i w_i \|Q\|_F^2 = 16H^2 \|Q\|_F^2 \right\} \right|. \end{aligned}$$

Let us now bound the first term in the RHS of (A.4) above. First, observe that  $x^\top Q x - \mathbb{E}_{\mathcal{N}(\mu_i, \Sigma_i)} x^\top Q x = (x - \mu_i)^\top Q (x - \mu_i) - \mathbb{E}_{\mathcal{N}(\mu_i, \Sigma_i)} (x - \mu_i)^\top Q (x - \mu_i) + 2(x - \mu_i)^\top Q \mu_i$ . Thus, using Fact 2.21 and Lemma 2.42, we have:

$$\left| \frac{Q}{2} \left\{ \sum_{i \leq k} w_i \mathbb{E}_{\mathcal{D}_i} \left( x^\top Q x - \mathbb{E}_{x \sim \mathcal{D}_i} x^\top Q x \right)^2 \right. \right. \quad (\text{A.5})$$

$$\leq 2 \sum_{i \leq k} w_i \mathbb{E}_{\mathcal{D}_i} \left( (x - \mu_i)^\top Q (x - \mu_i) - \mathbb{E}_{x \sim \mathcal{D}_i} (x - \mu_i)^\top Q (x - \mu_i) \right)^2 + 8 \sum_{i \leq k} w_i \mathbb{E}_{\mathcal{D}_i} \left( (x - \mu_i)^\top Q \mu_i \right)^2 \quad (\text{A.6})$$

$$\leq 6 \sum_i w_i C \left\| \Sigma_i^{1/2} Q \Sigma_i^{1/2} \right\|_F^2 + 8 \sum_{i \leq k} w_i \mathbb{E}_{\mathcal{D}_i} \left( (x - \mu_i)^\top Q \mu_i \right)^2. \quad (\text{A.7})$$

For the first term, note that  $\|\Sigma_i\|_2 \leq 1 + \|\Sigma_i - I\|_F \leq 1 + H$ . Thus,  $\left\| \Sigma_i^{1/2} \right\|_2 \leq \sqrt{1 + H}$ . Thus, we have that  $\Sigma_i^{1/2} \leq I + (\Sigma_i^{1/2} - I) \leq \sqrt{1 + H} I$ . Using Lemma 2.25 with  $A = (1 + H)^{-1/2} \Sigma_i^{1/2}$  and  $B = Q \Sigma_i^{1/2}$ , we have:  $\left| \frac{Q}{2} \left\{ \left\| \Sigma_i^{1/2} Q \Sigma_i^{1/2} \right\|_F^2 \leq (1 + H) \left\| Q \Sigma_i^{1/2} \right\|_F^2 \right\} \right|$ . By another application of Lemma 2.25, we have:

$$\left| \frac{Q}{2} \left\{ \left\| Q \Sigma_i^{1/2} \right\|_F^2 \leq (1 + H) \|Q\|_F^2 \right\} \right|. \quad \text{Thus, altogether, we have: } \left| \frac{Q}{2} \left\{ \left\| \Sigma_i^{1/2} Q \Sigma_i^{1/2} \right\|_F^2 \leq (1 + H)^2 \|Q\|_F^2 \right\} \right|.$$

Using our assumption that  $1 < H$ , we thus have:

$$\left| \frac{Q}{2} \left\{ \sum_i w_i C \left\| \Sigma_i^{1/2} Q \Sigma_i^{1/2} \right\|_F^2 \leq C(1 + H)^2 \|Q\|_F^2 \leq 4CH^2 \|Q\|_F^2 \right\} \right|.$$

For the second term, first observe that the following equality of quadratic polynomials in indeterminate  $Q$ :  $\left( (x - \mu_i)^\top Q \mu_i \right)^2 = \left( (\Sigma_i^{\dagger/2} (x - \mu_i))^\top \Sigma_i^{1/2} Q \mu_i \right)^2$ . Thus,  $\mathbb{E}_{x \sim \mathcal{D}_i} \left( (x - \mu_i)^\top Q \mu_i \right)^2 = \left\| \Sigma_i^{1/2} Q \mu_i \right\|_2^2$ . Next, by the SoS Cauchy-Schwarz inequality (Fact 2.20), we have that

$$\left| \frac{Q}{2} \left\{ \left\| \Sigma_i^{1/2} Q \mu_i \right\|_2^2 = \text{tr}(\mu_i \mu_i^\top Q \Sigma_i Q) \leq H \text{tr}(Q \Sigma_i Q) = H \left\| \Sigma_i^{1/2} Q \right\|_F^2 \right\} \right|.$$

Applying Lemma 2.25 with the observation above that  $\Sigma_i^{1/2} \leq (1+H)^{1/2}I$  yields:  $\left\{ \left\| \Sigma_i^{1/2} Q \right\|_F^2 \leq (1+H) \|Q\|_F^2 \right\}$ . Thus, altogether, we obtain:  $\left| \frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}_i} (x - \mu_i)^\top Q \mu_i \right\}^2 \leq H(1+H)^2 \|Q\|_F^2 \leq 4H^3 \|Q\|_F^2 \right.$ . We thus have:

$$\left. \left\{ \sum_{i \leq k} w_i \mathbb{E}_{\mathcal{D}_i} ((x - \mu_i)^\top Q \mu_i)^2 \leq 4H^3 \|Q\|_F^2 \right\} \right\}.$$

Plugging in these bounds into (A.7) completes the proof.  $\square$

As an immediate corollary of Lemma 2.39 and Lemma 2.43, we obtain:

**Lemma A.11** (Variance of Degree-2 Polynomials of Mixtures of Gaussians, Lemma 2.44 restated). *Let  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture of Gaussians with  $w_i \geq \alpha$ , mean  $\mu = \sum_i w_i \mu_i$  and covariance  $\Sigma = \sum_i w_i ((\mu_i - \mu)(\mu_i - \mu)^\top + \Sigma_i)$ . Suppose that for some  $H > 1$ ,  $\left\| \Sigma^{1/2} (\Sigma_i - I) \Sigma^{1/2} \right\|_F \leq H$  for every  $1 \leq i \leq k$ . Let  $Q$  be a symmetric  $d \times d$  matrix-valued indeterminate. Then for  $H' = \max\{H, 1/\alpha\}$ ,*

$$\left| \frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{M}} \left[ \left( x^\top Q x - \mathbb{E}_{x \sim \mathcal{M}} [x^\top Q x] \right)^2 \right] \right\} \leq 100H'^2 \left\| \Sigma^{1/2} Q \Sigma^{1/2} \right\|_F^2 \right\}.$$

*Proof.* Let  $\Sigma = U\Lambda U^\top$  be the covariance of the mixture  $\mathcal{M}$  along with its eigendecomposition. We want to apply Lemma 2.43 and Lemma 2.39 with the linear transformation  $x \rightarrow Ax$  for  $A = \Lambda^{1/2} U^\top$ . For this, we need to check that the conditions of the Lemma 2.43 are met after this linear transformation. The new component covariance is  $\Sigma'_i = A \Sigma_i A^\top$  and the hypothesis implies that they are within  $H$  in Frobenius distance of the new mixture covariance  $I' = A \Sigma A^\top$  ( $I$  in the range space of  $\Sigma$ ). The new means of the components after the linear transformation are  $\mu'_i = A \mu_i$  and the new mixture mean is  $\mu' = A \mu$ . Thus, noting that  $I' = \sum_i w_i (\mu'_i - \mu')(\mu'_i - \mu')^\top + \sum_i w_i \Sigma'_i$ , and since each of the terms in the RHS of the preceding equality are PSD, we must have that  $I' \geq \sum_i w_i (\mu'_i - \mu')(\mu'_i - \mu')^\top$  for every  $i$ . Thus,  $1 = \|I'\|_2 \geq w_i \|(\mu'_i - \mu')(\mu'_i - \mu')^\top\|_2 = \|\mu'_i - \mu'\|_2^2$ . Rearranging yields that  $\|\mu'_i - \mu'\|_2^2 \leq 1/w_i \leq 1/\alpha$ . Thus, we can now apply Lemma 2.43 to the linearly transformed mixture and the conclusion follows.  $\square$

#### A.4 Omitted Proofs from Section 2.4

**Lemma A.12** (Lemma 2.49 restated). *If  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$  with parameters  $(\gamma, t)$ , then if  $w_i \geq \gamma$  for all  $i \in [k]$ , and if for some  $B \geq 0$  we have that  $\|\mu_i\|_2^2, \|\Sigma_i\|_{\text{op}} \leq B$  for all  $i \in [k]$ , then for all  $m \leq t$ , we have that:*

$$\left\| \mathbb{E}_{x \in_u X} [x^{\otimes m}] - \mathbb{E}_{x \sim \mathcal{M}} [x^{\otimes m}] \right\|_F^2 \leq \gamma^2 m^{O(m)} B^m d^m.$$

*Proof.* We begin by noting that for any symmetric  $m$ -tensor  $A$  we have that  $\|A\|_F^2 \leq m^{O(m)} (\mathbb{E}_{v \sim \mathcal{N}(0, I)} [\langle A, v^{\otimes m} \rangle^2])$ . This is because, in the notation of [DKS18], the squared expectation of  $\langle A, v^{\otimes m} \rangle$  is  $\mathbb{E}[\text{Hom}_A(v)^2] \geq m! \mathbb{E}[h_A(v)^2] = m! \|A\|_F^2$ , where the first inequality holds

because  $\sqrt{m!}h_A(v)$  is the degree- $m$  harmonic part of  $\text{Hom}_A(v)$ , and the equality is by Claim 3.22. Therefore, to prove the lemma, it suffices to bound

$$\begin{aligned}
& \mathbb{E}_{v \sim \mathcal{N}(0, I)} \left[ \left( \mathbb{E}_{x \in \cup X} [(v \cdot x)^m] - \mathbb{E}_{x \sim \mathcal{M}} [(v \cdot x)^m] \right)^2 \right] \\
&= \mathbb{E}_{v \sim \mathcal{N}(0, I)} \left[ \left( \sum_{i=1}^k \left( \frac{1}{n} \sum_{x \in X_i} (v \cdot x)^m - w_i \mathbb{E}_{x \sim \mathcal{N}(\mu_i, \Sigma_i)} (v \cdot x)^m \right) \right)^2 \right] \\
&= \mathbb{E}_{v \sim \mathcal{N}(0, I)} \left[ \left( \sum_{i=1}^k \sum_{j=0}^m \binom{m}{j} \mu_i^{m-j} \left( \frac{1}{n} \sum_{x \in X_i} (v \cdot (x - \mu_i))^j - w_i \mathbb{E}_{x \sim \mathcal{N}(\mu_i, \Sigma_i)} (v \cdot (x - \mu_i))^j \right) \right)^2 \right] \\
&\leq \mathbb{E}_{v \sim \mathcal{N}(0, I)} \left[ \left( \sum_{i=1}^k \sum_{j=0}^m \binom{m}{j} |v \cdot \mu_i|^{m-j} \left( w_i \gamma m! (v^T \Sigma v)^{j/2} \right) \right)^2 \right] \\
&\leq \gamma^2 m^{O(m)} \mathbb{E}_{v \sim \mathcal{N}(0, I)} \left[ \sum_{i=1}^k \left( w_i (|v \cdot \mu_i| + (v^T \Sigma v)^{1/2})^{2m} \right) \right] \\
&\leq \gamma^2 m^{O(m)} \mathbb{E}_{v \sim \mathcal{N}(0, I)} [2B^m \|v\|_2^{2m}] \\
&\leq \gamma^2 m^{O(m)} B^m d^m.
\end{aligned}$$

This completes the proof.  $\square$

**Lemma A.13** (Lemma 2.50 restated). *Let  $\mathcal{M} = \sum_i w_i \mathcal{N}(\mu_i, \Sigma_i)$ . Let  $S \subset [k]$  with  $\sum_{i \in S} w_i = w$ , and let  $\mathcal{M}' = \sum_{i \in S} (w_i/w) \mathcal{N}(\mu_i, \Sigma_i)$ . Then if  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$  for some  $\gamma < 1/(2k)$  with the corresponding partition being  $X = X_1 \cup X_2 \cup \dots \cup X_k$ , then  $X' = \bigcup_{i \in S} X_i$  satisfies Condition 2.47 with respect to  $\mathcal{M}'$  with parameters  $(O(k\gamma/w), t)$ .*

*Proof.* After noting that  $|X'| = w|X|(1 + O(k\gamma/w))$ , the rest follows straightforwardly from the definitions using the partition  $X' = \bigcup_{i \in S} X_i$ .  $\square$

**Lemma A.14** (Lemma 2.51 restated). *Let  $\mathcal{M} = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, \Sigma_i)$  and let  $n$  be an integer at least  $kt^{Ct} d^t / \gamma^3$ , for a sufficiently large universal constant  $C > 0$ , some  $\gamma > 0$ , and some  $t \in \mathbb{N}$ . If  $X$  consists of  $n$  i.i.d. samples from  $\mathcal{M}$ , then  $X$  satisfies Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$  with high probability.*

*Proof.* We will show that Condition 2.47 holds with high probability using that partition where  $X_i$  is the set of samples drawn from the  $i$ -th component of  $\mathcal{M}$ . Note that the second part of Condition 2.47 holds with high probability, so long as  $n$  is a sufficiently large multiple of  $d/\gamma^2$  by the VC-Theorem [DL01]. In particular, if we think of samples as being drawn from  $\mathbb{R}^d \times [k]$ , where the second coordinate denotes the component that the sample was drawn from, the second part of Condition 2.47 says that the empirical probability of any event  $H \times \{i\}$  is correct to within additive error  $\gamma$ . It is easy to see and well-known that the class of such events has VC-dimension  $O(d)$ , from which the desired bound follows.

For the first part of Condition 2.47, we claim that it holds with high probability so long as  $n \geq kt^{Ct}d^t/\gamma^3$ . To prove this, we show it separately for each  $i$  with  $w_i \geq \gamma$  (as otherwise there is nothing to prove) and take a union bound. As Condition 2.47 is invariant under affine transformations, we may perform an invertible affine transformation so that  $\mu_i = 0$  and  $\Sigma_i$  is the projection onto the first  $d'$  coordinates, for some  $d'$ . It is clear that only the first  $d'$  coordinates of any element of  $X_i$  will be non-zero. We claim that the first part of our condition will follow for a given  $m$ , so long as  $\|X_i\|/n - w_i \leq \gamma w_i$  (which holds with high probability if  $n \gg \log(k)/\gamma^3$ ), and

$$\left\| \mathbb{E}_{x \in_u X_i} [x^{\otimes m}] - \mathbb{E}_{x \sim \mathcal{N}(0, I_{d'})} [x^{\otimes m}] \right\|_F^2 \leq \gamma^2, \quad (\text{A.8})$$

as  $\frac{1}{n} \sum_{x \in X_i} \langle v, x - \mu_i \rangle^m = w_i(1 \pm \gamma) \langle \mathbb{E}_{x \in_u X_i} [x^{\otimes m}], v^{\otimes m} \rangle$ . It is easy to see that each entry of the tensor on the left hand side of Equation (A.8) has mean 0 and variance  $m^{O(m)}/|X_i|$ , and thus the expected size of the left hand side is  $m^{O(m)}d^m/|X_i|$ . Then, when  $n \geq k^{Ck}d^{4k}/\gamma^3$  for a sufficiently large constant  $C$ , all parts of our condition hold with high probability. This completes the proof.  $\square$

## A.5 Omitted Proofs from Section 6

**Lemma A.15** (Frobenius Distance to TV Distance, Lemma 6.5, restated). *Suppose  $\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)$  are Gaussians with  $\|\mu_1 - \mu_2\|_2 \leq \delta$  and  $\|\Sigma_1 - \Sigma_2\|_F \leq \delta$ . If the eigenvalues of  $\Sigma_1$  and  $\Sigma_2$  are at least  $\lambda > 0$ , then*

$$d_{\text{TV}}(\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)) = O(\delta/\lambda).$$

*Proof.* By Fact 2.1, we have

$$d_{\text{TV}}(\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)) = O\left(\left((\mu_1 - \mu_2)^\top \Sigma_1^{-1}(\mu_1 - \mu_2)\right)^{1/2} + \|\Sigma_1^{-1/2} \Sigma_2 \Sigma_1^{-1/2} - I\|_F\right).$$

Then the first term is  $\langle \mu_1 - \mu_2, \Sigma_1^{-1}(\mu_1 - \mu_2) \rangle^{1/2} \leq (\|\Sigma_1^{-1}\|_{\text{op}} \|\mu_1 - \mu_2\|_2^2)^{1/2} \leq \delta/\sqrt{\lambda}$ . The second term is

$$\begin{aligned} \|\Sigma_1^{-1/2} \Sigma_2 \Sigma_1^{-1/2} - I\|_F^2 &= \left\| \Sigma_1^{-1/2} (\Sigma_1 - \Sigma_2) \Sigma_1^{-1/2} \right\|_F^2 \\ &= \text{tr} \left( \left( \Sigma_1^{-1/2} (\Sigma_1 - \Sigma_2) \Sigma_1^{-1/2} \right)^2 \right) \\ &\leq \text{tr} \left( (\Sigma_1 - \Sigma_2)^2 \right) (1/\lambda)^2 \\ &\leq (\delta/\lambda)^2. \end{aligned}$$

Thus,

$$d_{\text{TV}}(\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)) = O(\delta/\sqrt{\lambda} + \delta/\lambda) = O(\delta/\lambda).$$

$\square$

**Lemma A.16** (Component Moments to Mixture Moments, Lemma 6.6 restated). *Let  $\mathcal{M} = \sum_{i \in [k]} w_i \mathcal{N}(\mu_i, \Sigma_i)$  be a  $k$ -mixture such that  $w_i \geq \alpha$ , for some  $0 < \alpha < 1$ , and  $\mathcal{M}$  has mean  $\mu$  and covariance*

$\Sigma$  and for all  $i \neq j \in [k]$ ,  $\|\Sigma^{+1/2}(\Sigma_i - \Sigma_j)\Sigma^{+1/2}\|_F \leq 1/\sqrt{\alpha}$ . Let  $X$  be a multiset of  $n$  samples satisfying Condition 2.47 with respect to  $\mathcal{M}$  with parameters  $(\gamma, t)$ , for  $0 < \gamma < (dk/\alpha)^{-ct}$ , for a sufficiently large constant  $c$ , and  $t \in \mathbb{N}$ . Let  $\mathcal{D}$  be the uniform distribution over  $X$ . Then,  $\mathcal{D}$  is  $2t$ -certifiably  $(c/\alpha)$ -hypercontractive and for  $d \times d$ -matrix-valued indeterminate  $Q$ ,  $\left| \frac{Q}{2} \left\{ \mathbb{E}_{\mathcal{M}}(x^\top Qx - \mathbb{E}_{\mathcal{M}} x^\top Qx)^2 \leq \mathcal{O}(1/\alpha) \|\Sigma^{1/2}Q\Sigma^{1/2}\|_F^2 \right\} \right.$

*Proof.* First, since  $\mathcal{M}$  is a  $k$ -mixture of Gaussians with minimum mixing weight  $w_{\min} \geq \alpha$ , it follows from Corollary 2.35 that  $\mathcal{M}$  is  $t$ -certifiably  $(4/\alpha)$  hypercontractive. Further, since  $X$  satisfies Condition 2.47 with parameters  $(\gamma, t)$ , it follows from Lemma 2.48 that the set  $X' = \{\Sigma^{+1/2}(x_i - \mu)\}_{x_i \in X}$  also satisfies Condition 2.47 with parameters  $(\gamma, t)$  w.r.t.  $\mathcal{M}' = \sum_{i \in [k]} w_i \mathcal{N}(\Sigma^{+1/2}(\mu_i - \mu), \Sigma^{+1/2}\Sigma_i\Sigma^{+1/2})$ . Since  $\|\Sigma^{+1/2}\Sigma_i\Sigma^{+1/2}\|_{\text{op}} \leq \mathcal{O}(1/\alpha)$ , it follows from Lemma 2.49 that for all  $m \leq t$ ,  $\|\mathbb{E}_{x \in X'}[x^{\otimes m}] - \mathbb{E}_{x \sim \mathcal{M}'}[x^{\otimes m}]\|_F^2 \leq \gamma^2 m^{\mathcal{O}(m)} d^m (1/\alpha)^m$ . Since  $\gamma < (dk/\alpha)^{-\mathcal{O}(t)}$ , it follows from Fact 2.45 that  $X$  is  $2t$ -certifiably  $(c/\alpha)$ -hypercontractive.

By assumption, for all  $i \neq j \in [k]$ , we have that  $\|\Sigma^{+1/2}(\Sigma_i - \Sigma_j)\Sigma^{+1/2}\|_F \leq 1/\sqrt{\alpha}$ . We can now apply Lemma 2.44 to obtain

$$\left| \frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{M}} \left[ \left( x^\top Qx - \mathbb{E}_{x \sim \mathcal{M}} [x^\top Qx] \right)^2 \right] \leq \mathcal{O}(1/\alpha) \|\Sigma^{1/2}Q\Sigma^{1/2}\|_F^2 \right\} \right.$$

Therefore, it follows from Fact 2.45 that since  $X$  satisfies Condition 2.47 with parameters  $(\gamma, t)$ , the uniform distribution  $\mathcal{D}_X$  on  $X$ ,  $\left| \frac{Q}{2} \left\{ \mathbb{E}_{x \sim \mathcal{D}_X} (x^\top Qx - \mathbb{E}_{x \sim \mathcal{D}_X} x^\top Qx)^2 \leq \mathcal{O}(1/\alpha) \|\Sigma^{1/2}Q\Sigma^{1/2}\|_F^2 \right\} \right.$   $\square$

## B Bit Complexity Analysis

Here we address numerical issues related to our computation. We begin with the assumption that the eigenvalues of our covariance matrices are bounded below.

**Lemma B.1.** *If  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  is a mixture of Gaussians  $G_i$  where each  $G_i$  has mean and covariance of norm at most  $2^b$  for some positive integer  $b$  and each  $G_i$  has covariance matrix whose eigenvalues are bounded below by some  $\lambda > 0$ . Let  $\mathcal{M}'$  be an  $\varepsilon$ -corruption of  $\mathcal{M}$  whose outputs are bounded by  $2^{\mathcal{O}(b)}$ . Let  $N$  be a sufficiently large polynomial in  $d^k/\varepsilon$  and let  $\eta$  be  $\lambda$  divided by a sufficiently large polynomial in  $2^b d/\varepsilon$  (where sufficiently large is degree  $\mathcal{O}(1)$ ). Then if our algorithm is given  $N$  i.i.d. samples from  $\mathcal{M}'$  with each of their coordinates rounded to a nearby multiple of  $\eta$  (by which we mean one of the two closest), then our algorithm runs in time  $\text{poly}(N, b, \log(1/\eta))$  and with high probability returns a list of mixtures of Gaussians  $X_i$  with at least one of the  $X_i$   $\text{poly}_k(\varepsilon)$ -close to  $\mathcal{M}$  in parameter distance.*

*Proof.* This follows from noting firstly that with high probability the any subset of the rounded samples will have moments  $\lambda/\text{poly}(d/\varepsilon)$ -close to their moments before rounding. This means that with high probability these rounded samples will satisfy Condition 2.47. This means that our algorithm satisfies the necessary correctness guarantees. Furthermore, given that our samples now all have bounded bit complexity, it is easy to see that the runtime of our algorithm is polynomial in  $N$  and the bit complexity.  $\square$



More generally, as long as the parameters of the components of our mixture can be expressed with bounded bit complexity, we can prove a similar result, without needing any lower bound on the covariances.

**Theorem B.2.** *Let  $\mathcal{M} = \sum_{i=1}^k w_i G_i$  be a mixture of Gaussians where the  $G_i$  are Gaussians whose means and covariance matrices can all be written with coefficients given by rational numbers with bit complexity at most  $b$  for some integer  $b$ . Let  $\mathcal{M}'$  be an  $\varepsilon$ -corruption of  $\mathcal{M}$  so that with probability 1 the returned points have size  $2^{O(b)}$ . Let  $N$  be a sufficiently large polynomial in  $d^k/\varepsilon$ . Then there exists an algorithm that given  $b$  bit-oracle access to these samples runs in time  $\text{poly}(N, b)$  and with high probability returns a mixture of Gaussians  $X$  so that  $d_{\text{TV}}(X, \mathcal{M}) < \text{poly}_k(\varepsilon)$ .*

*Proof.* We begin by showing that we can find a list of hypotheses at least one of which is close. It is then straightforward to show that we can run a tournament over these hypotheses to find a specific one that works. We also assume for simplicity that each  $w_i$  is at least  $3\varepsilon$ .

We begin by setting  $\lambda$  to be  $2^{-b \cdot d^{kC}}$  for a sufficiently large constant  $C$ . By adding each sample to a random sample from  $N(0, \lambda I)$ , we can produce samples from  $\tilde{\mathcal{M}}'$ , and  $\varepsilon$ -corruption of  $\tilde{\mathcal{M}} = \sum_{i=1}^k w_i \tilde{G}_i$  where  $\tilde{G}_i$  is the convolution of  $G_i$  with  $N(0, \lambda I)$ . Note that  $\tilde{G}_i$  is a Gaussian whose covariance has eigenvalues at least  $\lambda$ . Furthermore, if the covariance matrix of  $G_i$  is non-singular, the smallest eigenvalue of the covariance matrix must be at least  $2^{O(b \cdot d)}$ , and therefore  $d_{\text{TV}}(G_i, \tilde{G}_i) < \varepsilon$ .

Since the eigenvalues of the components of  $\tilde{\mathcal{M}}$  are bounded below, we can apply Lemma B.1 to our samples from  $\tilde{\mathcal{M}}'$  rounded to an appropriate accuracy  $\eta$ , and in  $\text{poly}(N, b)$ -time obtain a list of hypothesis mixtures at least one of which is (with high probability) close to  $\tilde{\mathcal{M}}$  in total variation distance.

If the covariances of all of the  $G_i$  with weights more than some sufficiently large  $\text{poly}_k(\varepsilon)$  are all non-singular, then one of these hypotheses will be close to  $\mathcal{M}$ . Otherwise, there must be some  $i$  for which  $w_i$  is relatively large and for which  $G_i$  has singular covariance matrix. In particular, there must be an integer vector  $v$  with bit complexity  $O(bd)$  in the kernel of the covariance matrix of  $G_i$ . The hypothesis mixture  $X$  that is close to  $\tilde{\mathcal{M}}$  in parameter distance must contain some component close to  $\tilde{G}_i$ . Since  $\tilde{G}_i$  has covariance matrix  $\tilde{\Sigma}_i = \lambda I + \Sigma_i$  where  $\Sigma_i$  is the covariance matrix of  $G_i$ . We note that  $\tilde{\Sigma}_i$  will have an eigenvalue of  $\lambda$  and that therefore, our close hypothesis will have an eigenvalue at most  $2\lambda$ .

If any of our returned hypotheses have any component with a covariance matrix  $\Sigma$  which has any eigenvalue less than  $2\lambda$ , we do the following. We consider the quadratic form on integer vectors  $v$  defined by

$$Q(v) = v^T \Sigma v + \sqrt{\lambda} |v|_2^2.$$

We note that if this  $\Sigma$  is close in parameter distance to a singular  $\tilde{\Sigma}_i$  where  $\Sigma_i$  had a null-vector  $v$  of norm  $2^{O(bd)}$ , then for that same value of  $v$  we will have that  $Q(v) < \lambda^{1/4}$ . Using the Lovász local lemma in [EL73], we can find a  $v$  so that  $Q(v)$  is within a  $2^{O(d)}$ -factor of the minimum possible value over all non-zero, integer vectors  $v$ . If for this  $v$ ,  $Q(v) > 2^{\Omega(d)} \lambda^{1/4}$ , we know that the hypothesis in question is not close to  $\tilde{\mathcal{M}}$  in parameter distance and can be ignored. On the other hand, any  $v$  with  $Q(v)$  this small must have  $|v| < 2^{O(d)} \lambda^{-1/2}$  and  $v^T \Sigma v < 2^{O(d)} \lambda^{1/4}$ . Note that the projection of  $v$  onto the  $\ker(\Sigma_i)^\perp$  is either zero or has magnitude at least  $2^{O(bd)}$ . In the latter case, it would need to be the

case that  $Q(v)$  is substantially larger. Thus, if such a hypothesis is close to  $\tilde{\mathcal{M}}$  in parameter distance, then  $v$  is in the kernel of some  $\Sigma_i$ .

If our algorithm finds some  $v$  for some hypothesis, we then compute  $v \cdot x$  to error  $\lambda$  for each of our samples  $x$ . If  $\mathcal{M}$  really has a component with  $v$  in the kernel of its covariance matrix, all of the  $x$ 's taken from this component will have  $v \cdot x$  the same. This means that at least a  $(3/2)\varepsilon$  fraction of our samples  $x$  will have  $v \cdot x$  within  $\lambda$  of each other. Note that if  $v$  is not in the kernel of any covariance matrix of any  $G_i$  than  $\mathbf{Var}[v \cdot G_i]$  will be at least  $2^{O(db)}$  for each  $i$ , and with high probability we will not find this many close samples.

To summarize, if our algorithm applies this procedure to every component of every hypothesis and does not find such a  $v$ , then it cannot be the case that  $\mathcal{M}$  contains any components of weight more than  $\text{poly}_k(\varepsilon)$  that are singular, and thus one of our original hypotheses must be close in total variational distance. We can then run a tournament to find a single one that is close. Otherwise, if we find such a  $v$  for which many points do have  $v \cdot x$  close by, then  $v$  must be a null vector of the covariance matrix of some  $G_i$ . Furthermore, all of the samples within  $\lambda$  of this common value of  $v \cdot x$ , with high probability are either errors or come from components contained in some lower dimensional subspace. We can determine what this subspace is by noting that it is defined by  $v \cdot x = q$  for some rational number  $q$  with bit-complexity at most  $O(bd)$  and using continued fractions on a good numerical approximation of  $q$  in order to determine its true value. Our algorithm can then recurse on the points in this subspace (a mixture of Gaussians in a lower dimensional space) and on the remaining points (which are from a mixture of fewer Gaussians), and return an appropriate mixture of the results.  $\square$