Structures of Spurious Local Minima in k-means

Wei Qian, Yuqian Zhang, Yudong Chen

Abstract—The k-means clustering problem concerns finding a partition of the data points into k clusters such that the total within-cluster squared distance is minimized. This optimization objective is non-convex, and not everywhere differentiable. In general, there exist spurious local solutions other than the global optimum. Moreover, the simplest and most popular algorithm for k-means, namely Lloyd's algorithm, generally converges to such spurious local solutions both in theory and in practice. In this paper, we investigate the structures of these spurious local solutions under a probabilistic generative model with kground truth clusters. As soon as k = 3, spurious local minima provably exist, even for well-separated clusters. One such local minimum puts two centers at one true cluster, and the third center in the middle of the other two true clusters. We prove that this is essentially the only type of spurious local minima under a separation condition. In particular, any local minimum solution only involves a configuration that puts multiple centers at a true cluster, and one center in the middle of multiple true clusters. Our results pertain to the k-means formulation for mixtures of Gaussians or bounded distributions, and hold in the over- and under-parametrization regimes where the number of centers in k-means may not equal to the number of true clusters. Our theoretical results corroborate existing empirical observations and provide justification for popular heuristics for k-means clustering.

Index Terms—k-means, Gaussian mixture, non-convex optimization, spurious local minima, structured minima

I. INTRODUCTION

k-means clustering is one of the most fundamental problems in unsupervised learning, with a wide range of applications in multiple fields including machine learning, image analysis, computer graphics and beyond; see the survey [I] and the references therein. The *k*-means problem can be formulated as follows: given *n* data points $x_1, \ldots, x_n \in \mathbb{R}^d$, find *k* centers $\beta = (\beta_1, \ldots, \beta_k) \in \mathbb{R}^{d \times k}$ such that the following sum of squared distances is minimized.

$$G_n(\boldsymbol{\beta}) := \sum_{i=1}^n \min_{j \in [k]} \|\boldsymbol{x}_i - \boldsymbol{\beta}_j\|^2,$$
(1)

W. Qian was with the School of Operations Research and Information Engineering, Cornell University, Ithaca, NY, 14853 (email: wq34@cornell.edu). Y. Chen is with the School of Operations Research and Information Engineer-

ing, Cornell University, Ithaca, NY, 14853 (email: yudong.chen@cornell.edu.

Y. Zhang is with the Department of Electrical and Computer Engineering, Rutgers University, New Brunswick, NJ, 08854 (email: yqz.zhang@rutgers.edu).

W. Qian and Y. Chen are partially supported by NSF grants 1657420 and 1704828 and CAREER Award CCF-2047910.

Copyright (c) 2017 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

¹Another common way of formulating the *k*-means problem involves finding a partition of the data points into *k* clusters such that the within-cluster sum of squared distance is minimized. This partition-based formulation is equivalent to the center-based formulation (Π) used in this paper, as we show in Appendix [A]

where $\|\cdot\|$ denotes the ℓ_2 norm. The k-means objective function (1) is non-convex: it involves the minimization of quadratic functions and is symmetric with respect to permutation of the indices of components of β . This optimization problem is NP-hard in general [2, 3, 4]. It has been observed that standard algorithms for k-means often converge to spurious local solutions of (1) that are not globally optimal [5, 6]. Moreover, these local minima of k-means are prevalent in practice [7, 8].

Recent theoretical work has made progress in understanding the k-means and related clustering problems with two clusters. In particular, if the data is generated from a balanced mixture of two identical and spherical Gaussians, the work in [9] [10] [11]effectively shows that there is no spurious local minima, and that greedy algorithms such as the Lloyd's algorithm and Expectation-Maximization (EM) are guaranteed to converge to a global minimizer from a random initialization. However, as soon as there are more than two clusters, non-trivial spurious local solutions do exist, even when the ground truth clusters are well-separated and balanced. Worst yet, these spurious local solutions may have objective values arbitrarily worse than the global optimum, and randomly-initialized greedy algorithms may provably converge to these local solutions with high probability [12].

Despite the above negative results, not all hope is lost. In this paper, we show that even with a general number of clusters, a lot can be said about the structural properties of these spurious local minima. In particular, under certain mixture models, we prove that *all* spurious local minima of k-means are well-behaved, in the sense that they possess the same type of structure that partially recovers the global minimum. We elaborate below.

A. Main Contributions

Consider the k-means problem under the following probabilistic generative model. Let $\beta_1^*, \ldots, \beta_k^* \in \mathbb{R}^d$ be k^* distinct unknown true cluster centers. For each $s \in [k^*]$, let f_s be the density of a distribution with mean β_s^* . Each data point $x \in \mathbb{R}^d$ is sampled independently from a mixture f of these k^* distributions $\{f_s\}_{s \in [k^*]}$, with the density

$$f(\boldsymbol{x}) = \frac{1}{k^*} \sum_{s=1}^{k^*} f_s(\boldsymbol{x}).$$
 (2)

If each f_s is a Gaussian distribution centered at β_s^* , the above distribution becomes the (balanced/equally-weighted) Gaussian Mixture Model (GMM). Under the generative model (2), we consider the following population version of the *k*-means objective function:

$$G(oldsymbol{eta}) = \int \min_{j \in [k]} \|oldsymbol{x} - oldsymbol{eta}_j\|^2 f(oldsymbol{x}) \mathrm{d}oldsymbol{x}$$

$$= \frac{1}{k^*} \sum_{s=1}^{k^*} \int \min_{j \in [k]} \| \boldsymbol{x} - \boldsymbol{\beta}_j \|^2 f_s(\boldsymbol{x}) \mathrm{d} \boldsymbol{x}.$$
(3)

The objective function above can be viewed as the infinitesample $(n \to \infty)$ limit of the empirical objective function in equation (1). Note that this population objective is also non-convex.

Below for simplicity we assume $k = k^*$, that is, the number of fitted centers in the k-means formulation is the same as that of the true mixture, though our theoretical results hold more generally.

a) Existence of spurious local minima.: Under general conditions, the ground truth centers $\beta^* = (\beta_1^*, \ldots, \beta_k^*) \in \mathbb{R}^{d \times k}$ and any permutation thereof are (close to) a global minimum of G; see Proposition []. However, there exist additional spurious local minima, even in the simple one-dimensional setting with k = 3 clusters and when the densities $\{f_s\}_{s \in [k]}$ have bounded and disjoint supports. In particular, we show that one spurious local minimum $\beta = (\beta_1, \beta_2, \beta_3)$ has the following configuration:

$$\beta_1 \approx \beta_2 \approx \beta_1^*, \quad \beta_3 \approx \frac{\beta_2^* + \beta_3^*}{2}.$$
 (4)

In words, this local solution uses two centers to fit one true cluster, and the third center to fit the other two true clusters. See Proposition [2] for details. A similar observation was made in [12] for the log-likelihood objective function of Gaussian mixtures.

b) Structures of spurious local minima.: The local solution in equation (4) involves disjoint many-fit-one and one-fitmany associations. As our main result, we show that this is essentially the *only* type of spurious local minima for k-means under a separation condition:

Theorem (Informal). For well-separated mixture models, any non-degenerate local minima $\beta = (\beta_1, \ldots, \beta_k)$ of G only involves the following configurations: (i) multiple centers $\{\beta_j\}$ lie near a true cluster β_s^* , and (ii) one center β_j lies near the mean of multiple true clusters $\{\beta_s^*\}$. Moreover, the configurations (i) and (ii) involve disjoint sets of β_j 's and β_s^* 's.

See Theorems 1 and 2 for the precise statement of this result. In words, viewing a solution β as an assignment of the centers $\{\beta_j\}$ for fitting the ground truth clusters, we show that a local minimum β can only involve many-fit-one associations (case (i) above) and one-fit-many associations (case (ii) above), and each fitted center β_j or true center β_s^* only participates in one of these associations. Any other solution β with many-fit-many associations (or other configurations) *cannot* be a local minimum.

We illustrate the above results in Figure [] under a twodimensional GMM with 4 components. The top panels show different candidate solutions of k-means. The ground-truth centers are the only global minimum, as in Panel 1a. Panel 1b shows a spurious local minimum, where the orange center fits two clusters, and the blue and purple centers fit one cluster. In Panel 1c, the blue and green centers together fit 3 clusters; in Panel 1d, the blue and purple centers together fit 2 clusters. These two solutions contain many-fit-many associations and are *not* local minima.

For further verification, we run the Lloyd's algorithm [13] with the above four solutions as the initial solution. The Lloyd's algorithm is an iterative greedy method that alternates between assigning each data point to its closest center and updating the centers to be the means of the new clusters. It can be viewed as a quasi-Newton algorithm applied to the objective function (1) with a specific choice of step size **[14]**. The bottom panels in Figure **1** show the trajectories of intermediate solutions of Lloyd's algorithm and the final solutions they converge to. When initialized at a global or local minimum, the algorithm stays at the initial solutions as expected (Panels 2a and 2b). In Panel 2c, the algorithm escapes from the initial solution, which is not a local minimum, and then converges to the spurious local minimum plotted in Panel 1b. In Panel 2d, the algorithm again escapes from the initial solution and converges to the globally optimal ground-truth solution plotted in Panel 1a.

To put the above results in context, we note that even the existence (or the lack thereof) of spurious population local minima in GMM, posted as an open problem in the Conference on Learning Theory (COLT) 2007 [15], was not rigorously resolved until recently [10, [11, [12]]. As mentioned, in general spurious local minima do exist [12], as demonstrated by an example similar to that in equation [4]. Our results above provide a positive message in the context of the k-means objective: all local minima *partially* recover the global minimum, in the sense that they identify some of true cluster centers and the means of the other true cluster centers. Again see Panel 2b in Figure [1] for an illustration.

On the technical side, note that the k-means objective function G is complicated due to the min operator in its definition (3), and in particular is not everywhere differentiable. Our analysis makes uses of a new technique based on analyzing the Voronoi sets and boundaries associated with the min operator as well as the directional derivatives of certain smooth upper bounds of G. Doing so allows us to leverage both the first and second order optimality conditions to deduce the structures of a local minimum. We hope that this technique is useful for studying other mixture and clustering problems.

B. Related Work

With a history of more than 50 years [13, 5], the *k*-means problem has found broad applications in computer science, astronomy, biology, social science and beyond. We refer to the papers [16, 1] for a comprehensive survey of the work on this problem.

Without additional assumptions on the data points, optimizing the k-means objective in (1) is NP-hard when the number of components k is fixed [2] or when the dimension d is fixed [3]. Even finding a $(1 + \epsilon)$ approximation with varying (k, d) is hard [4]. Progress has been made on designing constant-ratio approximation algorithms; see, e.g., the results in [17, 18] among many others.

Lloyd's algorithm [13], often called *the k*-means algorithm, is arguably the most popular method for the *k*-means problem.



Figure 1. Top panels: Local minima and non-minima in GMM with 4 components. Solutions with many-fit-many configurations are not local minima. Bottom panels: Trajectories of greedy algorithm when initialized at different solutions. The colored circles correspond to an initial configuration of β . Running the Lloyd's k-means algorithm from this initialization converges to a solution denoted by colored squares. The black lines correspond to the trajectory of the intermediate iterates. The algorithm escapes from non-minima and converges to a global or local minimum.

In general, Lloyd's algorithm is only guaranteed to converge to a local minimum and is sensitive to initialization [19]. Moreover, it may take exponentially many steps to converge in the worst case [20, 21]. Under certain probabilistic assumptions of the data, several theoretical guarantees have been established for the Lloyd's algorithm [9, 22, 23]. In particular, the work in [22] shows that under a separation condition of the clusters, Lloyd's algorithm can learn a mixture of Gaussians as well as a mixture of bounded distributions, after (i) a dimension reduction step via PCA and (ii) an initialization that has a constant-factor approximation guarantee. There is also substantial work on designing provably efficient initialization schemes for Lloyd's algorithm [24, 25]. Particularly relevant to us is the work in [26], which considers over-parametrized k-means/EM (which fits k^* clusters using $k > k^*$ centers) equipped with extra pruning steps. Interestingly, the fitted centers they try to prune correspond to, in our language, many-fit-one associations (and sometimes almost-empty associations as well; see our main theorems). As Lloyd's algorithm finds local minima of k-means, our results can be used to characterize the structural properties of the output of Lloyd's. We emphasize that our results are in fact more general, applicable to the general k-means objective function (with or without over-parametrization) and hence not tied to a specific algorithm.

A recent line of work also considers convex relaxation methods for the *k*-means problem based on linear or semidefinite programming [27, 28, 29]. Theoretical guarantees have been established on when the solution of the convex program coincides with (or approximates) the global minimum of *k*-means [30, 31, 32, 33].

The *k*-means objective function can be viewed as a "hard" or limit version of the negative log-likelihood function for the Gaussian Mixture Model (GMM); see Section [II-B]. As such, our results are related to recent theoretical work on the Expectation-Maximization (EM) algorithm [34], which is a local/greedy algorithm for optimizing the likelihood function.

Positive results have been obtained on provable convergence of EM under GMM with k = 2 components [35], [10], [11], [36], [37]]. In particular, these results show that the negative log-likelihood function has no spurious local minima for a balanced mixture of two Gaussians with the same covariance matrix. However, in more general mixture models, it has been proved that spurious local minima do exist with high probability. Examples include a mixture of $k \ge 3$ equally weighted components [12], and a mixture of k = 2 unequally weighted components with known mixing weights [38].

II. PROBLEM SETUP

In this section, we introduce the statistical models for our main results. Recall the mixture model in equation (2) with k^* true clusters. We consider two concrete instantiations of this model.

A. Ball Mixture

The first instantiation is a mixture of uniform distributions on k^* disjoint balls. For each $u \in \mathbb{R}^d$, let $\mathbb{B}_u(r)$ denote the Euclidean ball centered at u with radius r. As the true centers $\{\beta_s^*\}_{s\in[k^*]}$ and the radius r are fixed throughout this paper, we use the shorthand $\mathbb{B}_s \equiv \mathbb{B}_{\beta_s^*}(r)$ for brevity. We assume that each data point x is sampled independently and uniformly from one of k^* disjoint balls centered at the true centers β_s^* ; that is, $x \sim \text{unif}(\mathbb{B}_s)$ with probability $\frac{1}{k^*}$.

This model, sometimes called the Stochastic Ball Model [30], is formally described below.

Definition 1 (Stochastic Ball Model). *The Stochastic Ball Model is the mixture* (2) *where each component has density*

$$f_s(\boldsymbol{x}) = \frac{1}{\operatorname{Vol}(\mathbb{B}_s)} \mathbb{1}_{\mathbb{B}_s}(\boldsymbol{x}), \quad s \in [k^*].$$

Here Vol(T) denotes the volume of a set $T \subseteq \mathbb{R}^d$ with respect to the Lebesgue measure, and $\mathbb{1}_T$ is the indicator function for the set T.

B. Gaussian Mixture

The second instantiation is the (spherical) Gaussian mixture model, where each data point \boldsymbol{x} is sampled independently from one of k^* Gaussian distributions whose means are the true centers $\{\boldsymbol{\beta}_s^*\}$; that is, $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\beta}_s^*, \sigma^2 \boldsymbol{I})$ with probability $\frac{1}{k^*}$. A formal description of GMM is given below.

Definition 2 (Gaussian Mixture Model). *The (spherical) Gaussian Mixture Model is the mixture* (2) *where each component has density*

$$f_s(\boldsymbol{x}) = rac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(-rac{\|\boldsymbol{x}-\boldsymbol{eta}^*_s\|^2}{2\sigma^2}
ight), \quad s\in[k^*].$$

We point out that the population negative likelihood function of GMM (with a posited variance parameter τ^2), namely²

$$L_{\tau}(\boldsymbol{\beta}) := -\int \log \left[\sum_{j \in [k]} \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{\beta}_j\|^2}{2\tau^2}\right)\right] f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x},$$

is closely related to the population k-means objective function G defined in equation (3). As the log-sum-exp function above is a form of soft maximum, L_{τ} can be viewed as a smooth approximation of G. Moreover, as $\tau \to 0$, we have $2\tau^2 L_{\tau}(\beta) \to G(\beta)$ for each β . In other words, the k-means objective function corresponds to the limit case of the GMM log-likelihood function, and hence results for one have immediate bearing upon the other.

C. Model Parameters

For both of the above models, we define the quantities

$$\Delta_{\max} := \max_{s \neq s'} \|\boldsymbol{\beta}_s^* - \boldsymbol{\beta}_{s'}^*\|,$$
$$\Delta_{\min} := \min_{s \neq s'} \|\boldsymbol{\beta}_s^* - \boldsymbol{\beta}_{s'}^*\|,$$

which are the maximum and minimum pairwise separations between the true centers. Accordingly, we introduce two quantities measuring the Signal-to-Noise Ratios (SNR) of the models. In particular, for the Stochastic Ball Model we define

$$\eta_{\max} := \frac{\Delta_{\max}}{r},$$
$$\eta_{\min} := \frac{\Delta_{\min}}{r},$$

which are the maximum and minimum separations normalized by the radius of the balls. For the Gaussian Mixture Model, we similarly define

$$\eta_{\max} := \frac{\Delta_{\max}}{\sigma \sqrt{\min(k + k^*, d)}},$$
$$\eta_{\min} := \frac{\Delta_{\min}}{\sigma \sqrt{\min(k + k^*, d)}}.$$

Note the $\sqrt{\min(k+k^*,d)}$ factor in the denominators above. This factor is the typical value of the norm of a random vector from a *d*-dimensional standard Gaussian distribution when projected to the $(k+k^*)$ -dimensional subspace spanned by the true and fitted centers $\{\beta_s^*\}_{s\in[k^*]}$ and $\{\beta_i\}_{i\in[k]}$. The above models are sometimes said to be *well-separated* if $\eta_{\min} = \Omega(1)$ [12]. Also note that the ratio $\frac{\eta_{\max}}{\eta_{\min}} \in [1, \infty)$ measures how evenly-spaced the true centers are. This ratio is close to 1 when the true centers are approximately equidistant to each other.

D. Voronoi sets

Each candidate solution $\beta = (\beta_1, \dots, \beta_k)$ of the *k*-means problem induces a *Voronoi diagram*, namely, a partition of the space \mathbb{R}^d based on proximity to the β_s 's. The Voronoi diagram plays a crucial role in understanding the *k*-means objective (3), which is defined by the quantity $\min_{j \in [k]} ||x - \beta_j||$, the distance of a point x to its closest center. Here we review some basic concepts related to Voronoi diagrams, which are useful for future development.

Given a set of k centers $\beta = (\beta_1, \dots, \beta_k) \in \mathbb{R}^{d \times k}$ in \mathbb{R}^d , let $\mathcal{V}_i(\beta)$ be the region consisting of points that are closer to β_i than to any other center β_j , $j \neq i$. Formally, for each $i \in [k]$ we define

$$\mathcal{V}_i(\boldsymbol{\beta}) := \left\{ \boldsymbol{x} \in \mathbb{R}^d : \|\boldsymbol{x} - \boldsymbol{\beta}_i\| \le \|\boldsymbol{x} - \boldsymbol{\beta}_j\|, \forall j \neq i \right\}.$$
(5)

We call each $\mathcal{V}_i(\beta)$ the Voronoi set associated with β_i . The Voronoi diagram of β is the collection of the Voronoi sets, that is, $\mathcal{V}(\beta) := \{\mathcal{V}_i(\beta) : i \in [k]\}$. Note that each Voronoi set is a polyhedron in \mathbb{R}^d with at most k - 1 facets, [3] as we can rewrite the definition in equation (5) as

$$\mathcal{V}_i(oldsymbol{eta}) = ig\{ oldsymbol{x} \in \mathbb{R}^d : 2\langle oldsymbol{eta}_j - oldsymbol{eta}_i, oldsymbol{x}
angle \leq \|oldsymbol{eta}_j\|^2 - \|oldsymbol{eta}_i\|^2, \ orall j
eq i, j \in [k] ig\}.$$

In addition, for each index pair (i, j) with $i \neq j$, we define the *Voronoi boundary* $\partial_{i,j}(\beta)$ as the intersection of the Voronoi sets associated with β_i and β_j ; that is,

$$\partial_{i,j}(oldsymbol{eta}) := \mathcal{V}_i(oldsymbol{eta}) \cap \mathcal{V}_i(oldsymbol{eta}) = ig\{oldsymbol{x} : \|oldsymbol{x} - oldsymbol{eta}_i\| = \|oldsymbol{x} - oldsymbol{eta}_j\|ig\}.$$

Note that $\partial_{i,j}(\beta)$ is the set of points with equal distance to β_i and β_j . If $\partial_{i,j}(\beta)$ has dimension d-1, we say that $\mathcal{V}_i(\beta)$ is *adjacent* to $\mathcal{V}_j(\beta)$, written as $\mathcal{V}_i(\beta) \sim \mathcal{V}_j(\beta)$. In this case, the two Voronoi sets $\mathcal{V}_i(\beta)$ and $\mathcal{V}_j(\beta)$ intersect at a common (full dimensional) facet of the two polyhedra. We use the notation $\partial(\beta) := \{\partial_{i,j}(\beta) : \mathcal{V}_i(\beta) \sim \mathcal{V}_j(\beta)\}$ to denote the collection of the Voronoi boundaries of adjacent Voronoi sets.

III. MAIN RESULTS

In this section, we present our main theoretical results on the structures of the local minima of the population k-means objective G defined in equation (3). In what follows, we use \mathbb{P} to denote the probability measure with respect to f, the density of the ground truth mixture. Similarly, for each $s \in [k^*]$, we use \mathbb{P}_s to denote the probability measure with respect to f_s , the density of the s-th component of the ground truth mixture.

²Here we ignore a constant additive term that is independent of the variable β .

³A facet is a (d-1) dimensional face of a *d*-dimensional polyhedron.



Figure 2. One-dimensional Stochastic Ball Model with radius r < 0.4 and ground truth cluster centers $\beta^* = (-2, 0, 2)$. The solution $\beta = (-2 - \frac{r}{2}, -2 + \frac{r}{2}, 1)$ is a spurious local minimum.

A. Stochastic Ball Model

Consider the Stochastic Ball Model in Definition \square We first state two simple results concerning the global and local minima of the *k*-means objective *G*. The first proposition, proved in Appendix \square states that the ground truth centers is the only global minimum of *G*.

Proposition 1 (Ground truth is global minimum). Under the Stochastic Ball Model with $k^* = k \ge 1$ and $\beta^* = (\beta_1^*, \beta_2^*, \ldots, \beta_k^*) \in \mathbb{R}^{d \times k}$. if $\eta_{\min} \ge 6\sqrt{k}$, then the true centers β^* (up to permutation of its k components) is the unique global minimum of G.

The next proposition, proved in Appendix A states that in general G has a spurious local minimum that is not a global minimum. An illustration is given in Figure 2.

Proposition 2 (Existence of spurious local minima). Consider the Stochastic Ball Model in one dimension with $k^* = k = 3$ and $\beta^* = (\beta_1^*, \beta_2^*, \beta_3^*) \in \mathbb{R}^{1\times 3}$ with $\beta_1^* = -2$, $\beta_2^* = 0$, $\beta_3^* = 2$, where each ground truth ball/interval has radius r. When r < 0.4 or equivalently $\eta_{\min} > 5$, the solution $\beta = (\beta_1, \beta_2, \beta_3) \in \mathbb{R}^{1\times 3}$ with $\beta_1 = -2 - \frac{r}{2}$, $\beta_2 = -2 + \frac{r}{2}$ and $\beta_3 = 1$ is a local minimum of G.

Proposition 2 is similar in spirit to the results in the work [12], which proves the existence of spurious local minima for the log-likelihood function of a Gaussian mixture. We note that their results are established using a limiting argument with $\eta_{\min} \rightarrow \infty$, whereas our result gives explicit values of β^* and β .

Conceptually, Proposition 1 shows that when $k^* = k$, the k-means objective G is a *statistically* sensible objective function for clustering, as its global minimum recovers the ground truth clustering. On the other hand, Proposition 2 highlights the *computational* difficulty of this optimization task, due to the existence of spurious local minima in the form of the configuration plotted in Figure 2

As the main result of this paper, we show that the above configuration is essentially the *only* local minimum. The following theorem, proved in Section VI holds in the general setting where the number of centers k in the k-means objective G is not necessarily equal to the number of clusters k^* in the true mixture.

Theorem 1 (Local minimum structures, Stochastic Ball Model). Under the Stochastic Ball Model, for each constant integer $c_0 \ge 1$ there exists a universal constant c > 0 for which the following holds. Assume that $\max(k^*, k) \le c_0$, $\eta_{\max} > 4c^2$ and $\eta_{\min} \ge 14c\sqrt{\eta_{\max}}$. If $\beta = (\beta_1, \ldots, \beta_k) \in \mathbb{R}^{d \times k}$ is a local minimum of G, then the ground truth centers and fitted centers can be partitioned as $[k^*] = \bigcup_{a=1}^m S_a^*$ and $[k] = \bigcup_{a=0}^m S_a$ respectively, such that for each $a \in [m]$, exactly one of the following holds:

• (many/one-fit-one association) $|S_a| \ge 1$ and $S_a^* = \{s\}$ for some $s \in [k]$; moreover,

$$\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_s^*\| \le \frac{10c}{\sqrt{\eta_{\max}}} \Delta_{\max} = 10c\sqrt{r\Delta_{\max}}, \quad \forall i \in S_a.$$

• (one-fit-many association) $S_a = \{i\}$ for some $i \in [k]$ and $|S_a^*| \ge 2$; moreover,

$$\left\|\boldsymbol{\beta}_{i} - \frac{1}{|S_{a}^{*}|} \sum_{s \in S_{a}^{*}} \boldsymbol{\beta}_{s}^{*}\right\| \leq \frac{14c}{\sqrt{\eta_{\max}}} \Delta_{\max} = 14c\sqrt{r\Delta_{\max}}.$$

In addition, for each $i \in S_0$, we have $\mathbb{P}(\mathcal{V}_i(\beta)) \leq \frac{c}{\sqrt{\eta_{\max}}}$ (almost-empty association).

Theorem \blacksquare states that all local minima have the same type of structure. In particular, if we view a candidate solution $\beta = (\beta_i)_{i=1}^k$ as configuring the centers β_i 's to fit the true clusters, then any local minimum β must be composed of only the following configurations:

- (i) many-fit-one: multiple β_i 's are close to the same ground truth center;
- (ii) one-fit-many: one β_i is close to the mean of several ground truth centers;
- (iii) almost-empty: a β_i is far (relatively to other β_j 's) from any ground truth center, in the sense that the Voronoi set of β_i is almost empty with a small measure.

In configurations (i) and (ii) above, being close means that the distance is order-wise smaller than the separation Δ_{\max} , by a factor of $1/\sqrt{\eta_{\max}}$; similarly, in configuration (iii), the measure being small means that it is on the order of $1/\sqrt{\eta_{\max}}$. Moreover, the configurations (i), (ii) and (iii) must involve disjoint sets of β_j s' and β_s^* s'. For concrete examples, recall Figure []; the ground truth solution in Panel 1a has 4 one-fit-one associations, whereas the spurious local minimum in Panel 1b consists of a two-fit-one, a one-fit-two and a one-fit-one association.

Put differently, Theorem 1 implies that if a solution β involves any configuration other than the three above, then β can be perturbed locally to strictly decreases its objective value. For example, the solutions in Panels 1c and 1d in Figure 1 use two centers to fit three and two true clusters, respectively. The objective value can be decreased by moving these two centers *away* from each other and towards different true clusters, as shown in Panels 2c and 2d. Our proof of Theorem 1 in fact makes use of this geometric idea, by testing the optimality conditions of the local minima along certain judiciously chosen directions.

It is instructive to specialize Theorem 1 to the limit case of a "point model", where $r \to 0$ or equivalently $\eta_{\max} \to \infty$; that is, each ground truth cluster s has a *point mass* at β_s^* . In this case, the three possibilities guaranteed in the theorem reduce to: (i) several β_i 's are located exactly at one true cluster β_s^* (many-fit-one); (ii) one center β_i is located at the mean of several true β_s^* 's (one-fit-many); (iii) for all the other β_i 's, their Voronoi sets do not contain any true clusters.

In the general setting with r > 0, Theorem [] guarantees that the above result for the point model still holds approximately, with an approximation error due to each true cluster having a mass spread around the true center. The three bounds in Theorem [] control the approximation errors with respect to cases (i)–(iii) in the point model above. These error bounds all scale with $1/\sqrt{\eta_{\text{max}}}$, which becomes smaller if the SNR η_{max} increases.

a) Tightness of the error bounds: The approximation errors above are unavoidable in general. We have already shown in Proposition 2 that there exists a local minimum $\beta = (\beta_1, \beta_2, \beta_3)$ where β_1 and β_2 are close but not exactly equal to β_1^* ; see Figure 2. Here the mass of the first true cluster \mathbb{B}_1 is equally split between the Voronoi sets of β_1 and β_2 , each of which lies at the corresponding center of mass (cf. Lemma 2), leading to a nonzero approximation error in the many-fit-one association. In addition, in Example 2 in Appendix OC we demonstrate another local minimum with a non-zero approximation error in the one-fit-many association.

Theorem [] gives upper bounds for these approximation errors, where the bounds take the form $1/\sqrt{\eta_{\text{max}}}$. The proof of the theorem in fact establishes a family of bounds that provide a trade-off between the errors for the three types of associations (see Theorem [4]). In particular, when $k^* = k$, for each number $\lambda \in (0, \frac{1}{2k^2r})$, one can derive the bounds

• many-fit-one – if $S_a^* = \{s\}$ and $|S_a| \ge 1$,

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{s}^{*}\| \leq O\left(\max\left(\frac{k}{\lambda}, kr + k^{2}\lambda r\Delta_{\max}\right)\right) \quad \forall i \in S_{a};$$

• one-fit-many – if $S_a = \{i\}$ and $|S_a^*| > 1$,

$$\left\|\boldsymbol{\beta}_{i} - \frac{1}{|S_{a}^{*}|} \sum_{s \in S_{a}^{*}} \boldsymbol{\beta}_{s}^{*}\right\| \leq O\left(kr + k^{2}\lambda r\Delta_{\max} + \frac{k}{\lambda}\right);$$

• almost-empty – $\mathbb{P}(\mathcal{V}_i(\boldsymbol{\beta})) \leq \lambda kr, \forall i \in S_0.$

where the partitions $[k^*] = \bigcup_{a=1}^m S_a^*$ and $[k] = \bigcup_{a=0}^m S_a$ may depend on λ . Taking $\lambda = \frac{c}{\sqrt{r\Delta_{\max}}} = \frac{c}{r\sqrt{\eta_{\max}}}$ gives the bounds in Theorem []. We are currently not sure whether these bounds are tight in general.

b) Separation condition: The result in Theorem 1 holds under the separation/SNR condition $\eta_{\min} = \Omega(1)$. Note that this kind of separation condition is needed even in recent work on more restrictive settings of mixture problems and EM/k-means, including those on $k^* = 2$ clusters [35, 39, 40, 31], on the local behavior of EM/k-means near the global optimum β^* [41, 42], and on the existence of spurious local minima [12] (the last work assumes $\eta_{\min} = \Omega(\sqrt{d})$). In our setting, we believe such a separation condition is in general necessary. Indeed, Example I in Appendix Oc shows that if η_{\min} is too small, then there exists a local minimum that involves a two-fit-three configuration and hence qualitatively violates the structural properties in Theorem I. Under our separation condition and additional assumptions, it may be possible to solve the mixture estimation problem using specialized algorithms with a good initialization (e.g., spectral initialization [43]). Nevertheless, as is the case in a recent line of work on the landscape of nonconvex problems [44, 45, 46, 47], our results are not tied to a specific algorithm, but rather concern the general geometry and structure of the k-means problem itself. We have not attempted to optimize the separation condition (particularly its scaling with k^*) and leave it as an important future problem whether

this condition can be improved. Indeed, we empirically observe that even with a much milder separation condition, such as the four clusters in Figure [I], the one-fit-many and many-fit-one configurations still hold for local minimum solution.

c) Over/under-parametrization: Theorem [] holds even when the number of centers k in the population k-means objective is different from the number of clusters k^* in the true mixture, including the over-parametrization regime $k > k^*$ and the under-parametrization regime $k < k^*$. As we discuss in greater details in Section [V] such flexibility has important algorithmic implications.

B. Gaussian Mixture Model

We next consider the Gaussian Mixture Model in Definition 2 The main difference between this model and the Stochastic Ball Model is that the Gaussian distribution has an unbounded support and thus the tails of the mixture components overlap with each other. Nevertheless, much of the results for the Ball Model can be extended to the Gaussian case. For example, one can establish results that are analogous to Propositions 1 and 2 regarding the global minima and the existence of spurious local minima. Here we focus on establishing an analogue of Theorem 1 which characterizes the structures of all local minima of the population k-means objective G.

Our main result is given in the following theorem, whose proof is given in Appendix A

Theorem 2 (Local minimum structures, Gaussian Mixture Model). Under the Gaussian Mixture Model, for each constant integer $c_0 \ge 1$ there exists a universal constant c > 0 for which the following holds. Assume that $\tilde{k} :=$ $\max(k, k^*) \le c_0$ and let $t \ge 1$ be any number satisfying $\varphi(t) := 2 \exp(-t^2 \min(d, \tilde{k})/8) < \frac{1}{4}$ and $7\tilde{k}\varphi(t) < 1$. Further assume that $\eta_{\max} \ge 16tc^2$ and $\eta_{\min} \ge 7c\sqrt{t\eta_{\max}} + 7\varphi(t)\eta_{\max}$. If $\beta = (\beta_1, \dots, \beta_k) \in \mathbb{R}^{d \times k}$ is a local minimum of G, then the ground truth centers and fitted centers can be partitioned as $[k^*] = \bigcup_{a=1}^m S_a^*$ and $[k] = \bigcup_{a=0}^m S_a$, respectively, such that for each $a \in [m]$, exactly one of the following holds:

• (many/one-fit-one association) $|S_a| \ge 1$ and $S_a^* = \{s\}$ for some $s \in [k]$; moreover,

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{s}^{*}\| \leq \Delta_{\max} \left\{ \frac{5c\sqrt{t}}{\sqrt{\eta_{\max}}} + 7c\varphi(t) \right\}, \quad \forall i \in S_{a}.$$
(6)

• (one-fit-many association) $S_a = \{i\}$ for some $i \in [k]$ and $|S_a^*| \ge 2$; moreover,

$$\left\|\boldsymbol{\beta}_{i} - \frac{1}{|S_{a}^{*}|} \sum_{s \in S_{a}^{*}} \boldsymbol{\beta}_{s}^{*}\right\| \leq \Delta_{\max} \left\{ \frac{7c\sqrt{t}}{\sqrt{\eta_{\max}}} + 7c\varphi(t) \right\}.$$
(7)

In addition, for each $i \in S_0$, we have $\mathbb{P}(\mathcal{V}_i(\boldsymbol{\beta})) \leq \frac{c\sqrt{t}}{\sqrt{\eta_{\max}}} + \varphi(t)$ (almost-empty association).

Theorem $\boxed{2}$ is qualitatively similar to Theorem $\boxed{1}$ and shows that the local minima in GMM have a similar type of structure. The only differences are that the separation condition in Theorem $\boxed{2}$ has an additional t factor, and that the bounds

for the three possibilities have an additional error term $\varphi(t)$ that decays exponentially in t^2 . The $\varphi(t)$ term reflects the influence of the exponential tail of a Gaussian distribution outside a ball of radius $t\sigma\sqrt{\min(d, k + k^*)}$. In fact, the proof of Theorem 2 proceeds by effectively reducing GMM to the Stochastic Ball Model, treating the bulk of the Gaussian as a bounded distribution and the tail as additional errors. The choice of t here determines what is viewed as the tail and hence controls the trade-off between the separation condition and the two terms in the error bounds. For a rough interpretation of the theorem, one could simply think of t as a numerical constant large enough so that $\varphi(t)$ is dominated by the other terms in the error bounds.

IV. IMPLICATIONS AND CONNECTIONS

The theorems in the last section provide *structural* results for the k-means objective. In this section, we discuss some *algorithmic* implications of these results and remark on their connections to the literature.

a) Algorithmic Implications: When $k^* = k$, our result implies that one can find the global minimum of k-means as long as the characteristic many-fit-one association for local minima can be avoided - in this case, one-fit-many association will also disappear as the number of true clusters and that of fitted centers are equal. This observation suggests that to avoid many-fit-one, one should initialize a greedy clustering algorithm without putting the fitted centers close to each other. Several popular heuristics for k-means implement precisely this idea. For instance, the celebrated k-means++ algorithm [24] is a version of the Lloyd's algorithm in which the initial centers are generated iteratively as follows: the first center is selected uniformly from the data points; after selecting j centers, one computes the minimal distance of each data point to these jcenters, and select a data point randomly as the (j + 1)-th center with probability proportional to the above distance. By design, this procedure tends to pick k initial centers that are far away from each other. Many other heuristics for k-means follow a similar spirit; see, e.g., the work in [48, 49, 50, 51].

On the other hand, our structural results also highlight the inherent combinatorial difficulty of the problem. In particular, when the number of clusters grows, there is a growing number of possible configurations with many-fit-one and one-fit-many associations. It then becomes easier to get trapped in one of the corresponding local minima. This is consistent with the systematic empirical study in [52], which observes that algorithms for k-means perform worse when there are more clusters.

b) Connection to Over-Parametrization: Our structural results in Theorem [] and [2] hold in the over-parametrization setting where $k > k^*$ centers are used to fit k^* ground truth clusters. Over-parametrization appears to be a promising approach for avoiding local minima. In particular, when k is sufficiently bigger than k^* , a random initial solution is likely to assign *at least* one center to each true cluster. In this case, one-fit-many association would be avoided. Running a greedy algorithm from this initial solution, one would expect that it converges to a solution with only many-fit-one and almostempty associations, which can then be pruned by inspecting

the pairwise distances of the fitted centers and the sizes of their Voronoi sets. The work in [26] implements this idea in the context of over-parametrized EM. In particular, after EM converges, they remove fitted centers with low mixing weights (corresponding to almost-empty association) and combine fitted centers that are close to each other (corresponding to many-fitone association).

In fact, the extensive empirical study in [53] shows that the above idea can be applied to other latent variable models, as these models often have a similar solution structure, i.e., some estimated latent variables having duplicated values or low prior probabilities.

V. PRELIMINARY PROPERTIES FOR THE k-means Objective

In this section, we derive several preliminary results on the analytical properties of the population k-means objective function G defined in equation (3), focusing on the Stochastic Ball Model. These properties are later used in the proofs of our main theorems.

When β has pairwise distinct components (i.e., $\beta_i \neq \beta_j, \forall i \neq j \in [k]$), it is often convenient to rewrite the function G using the notation of Voronoi sets:

$$G(\boldsymbol{\beta}) = \sum_{i=1}^{k} \int_{\mathcal{V}_{i}(\boldsymbol{\beta})} \|\boldsymbol{x} - \boldsymbol{\beta}_{i}\|^{2} f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}.$$
 (8)

We can see that G depends on β in a complicated way through both $||\boldsymbol{x} - \beta_i||^2$ and $\mathcal{V}_i(\beta)$. As shall become clear later, the dependence through the squared distance $||\boldsymbol{x} - \beta_i||^2$ determines the first-order condition for local optimality for G; on the other hand, understanding second-order conditions requires us to study the behaviors of the Voronoi sets $\mathcal{V}_i(\beta)$ under small perturbation of β . To deal with this complication, our main strategy is to understand the directional behaviors of G along certain (judiciously chosen) directions, and to construct upper bounds on G that are easier to work with.

A. Directional Behaviors of G

Throughout the remainder of this section, we fix a candidate solution $\boldsymbol{\beta} = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k) \in \mathbb{R}^{d \times k}$. For a given direction $\boldsymbol{v} = (\boldsymbol{v}_1, \boldsymbol{v}_2 \dots, \boldsymbol{v}_k) \in \mathbb{R}^{d \times k}$, we are interested in how the objective $G(\boldsymbol{\beta})$ changes after we perturb $\boldsymbol{\beta}$ to $\boldsymbol{\beta} + t\boldsymbol{v}$. Restricting the function G to the direction \boldsymbol{v} , we define the directional objective function

$$H^{\boldsymbol{v}}(t) := G(\boldsymbol{\beta} + t\boldsymbol{v}), \quad t \in \mathbb{R}.$$

Note that β is a local minimum of G if and only if **0** is local minimum of H^{v} for all v.

The functions G and H^{v} are not everywhere differentiable, as they involve the minimum of quadratic functions. However, they are differentiable almost everywhere. In particular, whenever β has pairwise distinct components, the directional derivative $\frac{d}{dt}H^{v}(0)$ is guaranteed to exist and admits a simple expression, as shown in the following lemma.

Lemma 1 (Directional derivative). Suppose that β satisfies $\beta_i \neq \beta_j$ whenever $i \neq j$. For any choice of direction v, the

directional derivative $\frac{d}{dt}H^{\boldsymbol{v}}(0)$ exists and has the following analytic formula:

$$\frac{\mathrm{d}}{\mathrm{d}t}H^{\boldsymbol{v}}(0) = -\sum_{i=1}^{k}\int_{\mathcal{V}_{i}(\boldsymbol{\beta})} 2\langle \boldsymbol{v}_{i}, \boldsymbol{x} - \boldsymbol{\beta}_{i}\rangle f(\boldsymbol{x})\mathrm{d}\boldsymbol{x}$$

The lemma follows from the Leibniz integral rule; we defer the proof to Appendix A. Note that the above expression only involves differentiating the integrand in the expression (8); the Voronoi sets $\mathcal{V}_i(\boldsymbol{\beta})$ remain unchanged when only the first-order derivative is concerned.

B. First-Order Necessary Condition for Local Optimality

Using the first-order derivative expression in Lemma 1 we can derive a necessary condition for β being a local minimum. In particular, the following lemma states that any local minimum (satisfying a non-degeneracy condition) must have pairwise distinct components, each of which must be the center of its Voronoi set.

Lemma 2 (Local minimum must be Voronoi centers). Suppose that β is a local minimum of G. Then for each pair $i \neq j$, we must have $\beta_i \neq \beta_j$ whenever $\mathcal{V}_i(\beta) \cup \mathcal{V}_j(\beta)$ having a positive measure (with respect to f). Moreover, for each β_i whose Voronoi set $\mathcal{V}_i(\beta)$ has a positive measure, β_i must be at the center of probability mass of the Voronoi set $\mathcal{V}_i(\boldsymbol{\beta})$; that is,

$$\boldsymbol{\beta}_{i} = \frac{\int_{\mathcal{V}_{i}(\boldsymbol{\beta})} \boldsymbol{x} f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}}{\int_{\mathcal{V}_{i}(\boldsymbol{\beta})} f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}}.$$
(9)

We prove Lemma 2 in Appendix A. The conclusion of the above lemma can be written equivalently in a more explicit way. In particular, for each pair $(i, s) \in [k] \times [k^*]$, let $m_{i,s}(\beta)$ and $c_{i,s}(\beta)$ denote the probability mass and the center of mass of the set $\mathcal{V}_i(\boldsymbol{\beta})$ with respect to f_s respectively:

$$egin{aligned} m_{i,s}(oldsymbol{eta}) &:= \int_{\mathcal{V}_i(oldsymbol{eta})} f_s(oldsymbol{x}) \mathrm{d}oldsymbol{x}, \ oldsymbol{c}_{i,s}(oldsymbol{eta}) &:= rac{\int_{\mathcal{V}_i(oldsymbol{eta})} oldsymbol{x} f_s(oldsymbol{x}) \mathrm{d}oldsymbol{x}}{m_{i,s}(oldsymbol{eta})}. \end{aligned}$$

Then equation (9) can be rewritten as

$$\boldsymbol{\beta}_{i} = \frac{\sum_{s=1}^{k^{*}} m_{i,s}(\boldsymbol{\beta}) \boldsymbol{c}_{i,s}(\boldsymbol{\beta})}{\sum_{s=1}^{k^{*}} m_{i,s}(\boldsymbol{\beta})}$$

C. Decomposition of H^{v}

Lemmas 1 and 2 provide a first-order characterization of the local minima of G. For a more precise characterization, we need to account for the change in the Voronoi sets $\mathcal{V}(\beta)$ and its boundaries $\partial(\beta)$ when perturbing β to $\beta + tv$. With t > 0considered arbitrarily small, we make two observations:

- 1) The Voronoi set boundaries $\partial(\beta + tv)$ change continuously with respect to t.
- 2) When β is perturbed by tv, the points swept by the boundaries $\partial(\beta + tv)$ change their association from one Voronoi set to another.

Formally, for each pair $(i, j) \in [k] \times [k]$ we define the set

$$\Delta_{i \to j}^{\boldsymbol{v}}(t) := \mathcal{V}_i(\boldsymbol{\beta}) \cap \mathcal{V}_j(\boldsymbol{\beta} + t\boldsymbol{v}),$$

which is the set of points that change association from the *i*-th fitted center to the *j*-th fitted center due to the perturbation tv. Being the intersection of two polyhedra, the set $\Delta_{i \to i}^{v}(t)$ is a also polyhedron. An illustration of $\Delta_{i \to j}^{v}(t)$ is provided in Figure 3

As previously shown in Lemma 2, any non-degenerate local minimum β must have distinct components, so the corresponding Voronoi sets are also distinct. The same holds for the perturbed solution $\beta + tv$ when t is sufficiently small. In this case, we can decompose the directional objective function $H^{\boldsymbol{v}}$ as follows:

$$H^{\boldsymbol{v}}(t) = \sum_{i=1}^{k} \int_{\mathcal{V}_{i}(\boldsymbol{\beta}+t\boldsymbol{v})} \|\boldsymbol{x}-\boldsymbol{\beta}_{i}-t\boldsymbol{v}_{i}\|^{2}f(\boldsymbol{x})d\boldsymbol{x}$$

$$= \underbrace{\sum_{i=1}^{k} \int_{\mathcal{V}_{i}(\boldsymbol{\beta})} \|\boldsymbol{x}-\boldsymbol{\beta}_{i}-t\boldsymbol{v}_{i}\|^{2}f(\boldsymbol{x})d\boldsymbol{x}}_{U^{\boldsymbol{v}}(t)} + \underbrace{\sum_{(i,j):i\neq j} \int_{\Delta_{i\rightarrow j}^{\boldsymbol{v}}(t)} \left(\|\boldsymbol{x}-\boldsymbol{\beta}_{j}-t\boldsymbol{v}_{j}\|^{2} - \|\boldsymbol{x}-\boldsymbol{\beta}_{i}-t\boldsymbol{v}_{i}\|^{2}\right)f(\boldsymbol{x})d\boldsymbol{x}}_{W^{\boldsymbol{v}}(t)}$$

$$(10)$$

Here $U^{\boldsymbol{v}}(t)$ and $W^{\boldsymbol{v}}(t)$ correspond to the change in the objective value from two different sources. In particular, $U^{\boldsymbol{v}}(t)$ is due to the change in the distance between the data points and the centers, and $W^{\boldsymbol{v}}(t)$ is due to the data points changing association with the Voronoi sets.

Remark 1. By definition of $\Delta_{i \to j}^{v}(t)$, for each $x \in \Delta_{i \to j}^{v}(t)$, the integrand $||x - \beta_j - tv_j||^2 - ||x - \beta_i - tv_i||^2$ in the definition of $W^{\boldsymbol{v}}(t)$ is non-positive.

Proof of equation (10). When the Voronoi sets $\mathcal{V}(\beta)$ are perturbed to $\mathcal{V}(\boldsymbol{\beta} + \overline{tv})$, each point \boldsymbol{x} in \mathbb{R}^d either remains associated with the *i*-th center for some *i*, or changes its association from the i-th center to the j-th center for some $j \neq i$. In the first case, since $x \in \mathcal{V}_i(\beta) \cap \mathcal{V}_i(\beta + tv)$, we see that the contribution from x to $H^{v}(t)$ appears in U^{v} . In the second case, since $x \in \Delta_{i \to j}^{v}(t) = \mathcal{V}_{i}(\beta) \cap \mathcal{V}_{j}(\beta + tv)$, we can write the contribution from x as

$$\begin{aligned} \|\boldsymbol{x} - \boldsymbol{\beta}_j - t\boldsymbol{v}_j\|^2 \\ = \|\boldsymbol{x} - \boldsymbol{\beta}_i - t\boldsymbol{v}_i\|^2 + (\|\boldsymbol{x} - \boldsymbol{\beta}_j - t\boldsymbol{v}_j\|^2 - \|\boldsymbol{x} - \boldsymbol{\beta}_i - t\boldsymbol{v}_i\|^2), \end{aligned}$$

which appears in both $U^{\boldsymbol{v}}$ and $W^{\boldsymbol{v}}$.

D. Smooth Upper Bounds of H^{v}

The expression (10) for H^{v} is quite complicated. To understand the local minima of H^{v} , we instead study a simpler, better-behaved upper bound function of H^{v} that preserves the local minima and is amenable to calculus tools. In particular, we make use of the following elementary lemma.

Lemma 3 (Smooth upper bound). Suppose that $h, \tilde{h} : \mathbb{R} \to \mathbb{R}$ are two continuous functions that satisfy $h \leq h$ and h(0) = $\tilde{h}(0)$. If 0 is a local minimum of h, then 0 is also a local

=



Figure 3. Illustration of the set $\Delta_{i \to j}^{v}(t)$. The red dots represent the original centers $\{\beta_i\}$ and the blue stars represent the perturbed centers $\{\beta'_i = \beta_i + tv_i\}$. The blue solid lines are the original Voronoi boundaries $\partial(\beta)$, and the red dashed lines are the perturbed Voronoi boundaries $\partial(\beta')$. Top panel: the moving direction $v = (v_1, v_2, v_3)$ satisfies $v_1 = v_2 = v_3$, in which case the Voronoi boundaries are shifted parallelly by tv_1 . Bottom panel: the moving direction satisfies $v_1 = -v_2 = -v_3$, in which case the boundary $\partial_{1,2}(\beta')$ rotates around the mid point $\frac{\beta_1 + \beta_2}{2}$, $\partial_{1,3}(\beta')$ rotates around the mid point $\frac{\beta_1 + \beta_3}{2}$, and $\partial_{2,3}(\beta')$ shifts parallelly in the direction of v_2 . Each colored region represents $\Delta_{i \to j}^{v}(t)$, the set of points that change the association from the *i*-th center to the *j*-th center.

minimum of \tilde{h} ; moreover, we have $\lim_{t\to 0} \frac{\tilde{h}(t) - \tilde{h}(0)}{t} = 0$ and $\lim_{t\to 0} \frac{\tilde{h}(t) - \tilde{h}(0)}{t^2} \ge 0$ whenever the limits exist.

Proof. Since 0 is a local minimum of h, we have $\tilde{h}(0) = h(0) \le h(t) \le \tilde{h}(t)$ for all t in a neighborhood of 0, so 0 is also a local minimum of \tilde{h} . The first-order optimality condition for 0 gives $\lim_{t\to 0} \frac{\tilde{h}(t) - \tilde{h}(0)}{t} = \tilde{h}'(0) = 0$. Moreover, we have $\tilde{h}(t) - \tilde{h}(0) \ge 0 \implies \frac{\tilde{h}(t) - \tilde{h}(0)}{t^2} \ge 0$ for all $t \ne 0$ in a neighborhood of 0, which implies that $\lim_{t\to 0} \frac{\tilde{h}(t) - \tilde{h}(0)}{t^2} \ge 0$.

With the above lemma, we can study the structure of each local minimum of H^{v} (and hence that of G) by exploiting the optimality conditions of a smooth upper bound of H^{v} that is tight at the minimum. Let us take a first step in constructing such an upper bound. In view of the non-positivity result in Remark [], we can obtain an upper bound of the function W^{v} defined in (10) by only considering those pairs (i, j) for which $\mathcal{V}_{i}(\beta) \sim \mathcal{V}_{i}(\beta)$ are adjacent:

$$\begin{split} & W^{\boldsymbol{v}}(t) \\ \leq & \sum_{\substack{(i,j):\\ \mathcal{V}_i(\boldsymbol{\beta}) \sim \mathcal{V}_j(\boldsymbol{\beta})}} \int_{\Delta_{i \to j}^{\boldsymbol{v}}(t)} \left(\|\boldsymbol{x} - \boldsymbol{\beta}_j - t\boldsymbol{v}_j\|^2 - \|\boldsymbol{x} - \boldsymbol{\beta}_i - t\boldsymbol{v}_i\|^2 \right) f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} \\ &= & \sum_{\substack{(i,j):\\ \mathcal{V}_i(\boldsymbol{\beta}) \sim \mathcal{V}_j(\boldsymbol{\beta})}} \frac{1}{k^*} \sum_{s=1}^{k^*} W_{i \to j,s}^{\boldsymbol{v}}(t) \end{split}$$

where the quantity $W_{i \to j,s}^{\boldsymbol{v}}$, a shorthand for

$$\int_{\Delta_{i
ightarrow j}^{oldsymbol{v}}(t)}(\|oldsymbol{x}-oldsymbol{eta}_j-toldsymbol{v}_j\|^2-\|oldsymbol{x}-oldsymbol{eta}_i-toldsymbol{v}_i\|^2)f_s(oldsymbol{x})\mathrm{d}oldsymbol{x},$$

represents the contribution from the points in the *s*-th true cluster that change association from the *i*-th center to the *j*-th center. Combining the above inequality with equation (10), we obtain the following upper bound

$$H^{\boldsymbol{v}}(t) \leq U^{\boldsymbol{v}}(t) + \sum_{\substack{(i,j):\\\mathcal{V}_i(\boldsymbol{\beta})\sim\mathcal{V}_j(\boldsymbol{\beta})}} \frac{1}{k^*} \sum_{s=1}^{k^*} W^{\boldsymbol{v}}_{i\to j,s}(t).$$
(11)

In the proofs of our main theorems, we build upon equation (11) to derive further smooth upper bounds of H^v .

VI. PROOF OF THEOREM 1

In this section, we prove our main result under the Stochastic Ball Model (Definition 1). We in fact establish a more general and quantitative version of Theorem 1 that holds for any k^* and k.

Theorem 3 (General version of Theorem 1). Under the Stochastic Ball Model, assume that $\eta_{\max} > 4c^2k^4$ and $\eta_{\min} \ge \sqrt{\eta_{\max}} \left(2ck^2 \left(1 + \frac{2k^*}{k} \right) + 4k^* \left(3 + \frac{2k^*}{k} \right) \right)$ for some universal constant $c \ge 3$. If $\beta = (\beta_1, \ldots, \beta_k) \in \mathbb{R}^{d \times k}$ is a local minimum of G, then the ground truth centers and fitted centers can be partitioned as $[k^*] = \bigcup_{a=1}^m S_a^*$ and $[k] = \bigcup_{a=0}^m S_a$ respectively, such that for each $a \in [m]$, exactly one of the following holds:

• (many/one-fit-one association) $|S_a| \ge 1$ and $S_a^* = \{s\}$ for some $s \in [k^*]$; moreover,

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{s}^{*}\| \leq \Delta_{\max} \frac{\left(4k^{*} + 2ck^{2}\right)\left(1 + \frac{2k^{*}}{k}\right)}{\sqrt{\eta_{\max}}}, \quad \forall i \in S_{a}$$

• (one-fit-many association) $S_a = \{i\}$ for some $i \in [k]$ and $|S_a^*| \ge 2$; moreover,

$$\left\|\boldsymbol{\beta}_{i} - \frac{1}{|S_{a}^{*}|} \sum_{s \in S_{a}^{*}} \boldsymbol{\beta}_{s}^{*}\right\| \leq \Delta_{\max} \frac{2ck^{2}\left(1 + \frac{2k^{*}}{k}\right) + 4k^{*}\left(3 + \frac{2k^{*}}{k}\right)}{\sqrt{\eta_{\max}}}.$$

In addition, for each $i \in S_0$, we have $\mathbb{P}(\mathcal{V}_i(\beta)) \leq \frac{ck}{\sqrt{\eta_{\max}}}$ (almost-empty association).

Given Theorem 3. Theorem 1 follows immediately. The rest of the section is devoted to proving Theorem 3.

Throughout the proof, let β be a fixed local minimum of the k-means objective G. Note that G is invariant under translation of the space and permutation of the true centers. Consequently, we may assume without loss of generality that $\beta_1^* = \mathbf{0}$ and $\max_{s \in [k]} \|\beta_s^*\| \le \Delta_{\max}$.

A. Notations

We use V_d to denote the volume of a unit ball in \mathbb{R}^d with respect to the Lebesgue measure. For a set $T \subset \mathbb{R}^d$, let $\operatorname{int}(T)$ denote its interior, and $\operatorname{ReVol}(T)$ denote the relative volume of T with respect to the Lebesgue measure on the affine hull of T, with the convention that $\operatorname{ReVol}(\emptyset) = 0$. For two vector $u, u' \in \mathbb{R}^d, \angle (u, u') := \operatorname{arccos}(\frac{u^\top u'}{\|u\|\| \|u'\|}) \in [0, \pi]$ is the angle between u and u'. For each tuple $(i, j, s) \in [k] \times [k] \times [k^*]$ with $i \neq j$, we use $\mathcal{L}_{i,j,s}(\beta)$ to denote the two-dimensional plane that contains β_i, β_j and β_s^* (if such a plane is not unique, we fix an arbitrary one). Since we are concerned with a fixed local minimum β , we sometimes suppress the dependency on β and write, for example, $\mathcal{V}_i \equiv \mathcal{V}_i(\beta), \ \partial_{i,j} \equiv \partial_{i,j}(\beta)$ and $\mathcal{L}_{i,j,s} \equiv \mathcal{L}_{i,j,s}(\beta)$.

B. Proof of Theorem 3

To prove Theorem 3 we establish an intermediate result as given in Theorem 4 which provides a family of bounds parametrized by $\lambda > 0$.

Theorem 4 (Family of bounds for ball model). Under the Stochastic Ball Model, let $\beta = (\beta_1, \ldots, \beta_k)$ be a local minimum of the k-means objective function G defined in (3) and $\lambda > 0$ be an arbitrary fixed number. For each $i, j \in [k]$ and $s \in [k^*]$, let $\rho_s(\partial_{i,j}) := \frac{1}{V_d r^d} \operatorname{ReVol}(\partial_{i,j} \cap \mathbb{B}_s)$. For each $i \in [k]$, define the sets

$$T_i := \{ s \in [k^*] : \mathcal{V}_i \cap \mathbb{B}_s \neq \emptyset \} \quad and$$
$$A_i := \{ s \in [k^*] : \beta_s^* \in \operatorname{int}(\mathcal{V}_i) \} \subseteq T_i.$$

Then the following is true for each $i \in [k]$.

1) If $\rho_s(\partial_{j,\ell}) > \lambda$ for some $s \in T_i$ and some pair (j, ℓ) , then

$$\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_s^*\| \le \frac{k^*}{\lambda} + 3r.$$

2) For each $s \in T_i$, if $\rho_s(\partial_{j,\ell}) \leq \lambda$ for all pair (j,ℓ) , then the following bounds hold:

$$\mathbb{P}_{s}(\mathcal{V}_{i}) \geq 1 - k^{2}\lambda r, \quad \text{if } s \in A_{i}, \\ \mathbb{P}_{s}(\mathcal{V}_{i}) \leq k\lambda r, \qquad \text{if } s \in T_{i} \setminus A_{i}$$

Furthermore, if $\rho_s(\partial_{j,\ell}) \leq \lambda$ for all $s \in T_i$ and all pair (j, ℓ) , then:

a) When $|A_i| = 0$, we have

$$\mathbb{P}(\mathcal{V}_i) \le k\lambda r.$$

b) When $|A_i| > 0$, we have

$$\begin{aligned} \|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}\| &\leq \frac{k^{*}r}{1 - k^{2}\lambda r} \\ &+ \frac{k^{*}(k^{2}\lambda r^{2})(1 + (k^{*}k - k^{2})\lambda r)}{(1 - k^{2}\lambda r)^{2}} \\ &+ \frac{(k + 2k^{*})k\lambda r}{1 - k^{2}\lambda r}\Delta_{\max}, \end{aligned}$$
where $\boldsymbol{b}_{i} := \frac{1}{1 - k^{2}\lambda r}\sum_{a \neq a} \boldsymbol{\beta}^{*}_{a}$

where $\mathbf{b}_i := \frac{1}{|A_i|} \sum_{s \in A_i} \boldsymbol{\beta}_s^{-}$. The proof of Theorem \mathbf{A} which lies at

The proof of Theorem $\frac{4}{4}$, which lies at the core of our analysis, is given in Section VI-C.

We now derive Theorem 3 from Theorem 4 Doing so involves several elementary though somewhat tedious steps of a combinatorial flavor. To this end, we fix $\lambda = \frac{c}{\sqrt{r\Delta_{\max}}} = \frac{c}{r\sqrt{\eta_{\max}}}$. Recall the assumption in the main theorem that $\eta_{\max} > 4c^2k^4$ for c > 3. This assumption implies that $k^2\lambda r < 0.5$. If $\rho_s(\partial_{i,j}) > \lambda$, we say that a true cluster \mathbb{B}_s encloses the Voronoi boundary $\partial_{i,j}$ with a large relative volume; otherwise, we say that \mathbb{B}_s encloses the Voronoi boundary $\partial_{i,j}$ with a small relative volume.

We first state two simple implications of Theorem $\frac{4}{4}$ used frequently in the subsequent proof.

- Observation 1. For each $i \in [k]$, there exists at most one $s \in T_i$ such that $\rho_s(\partial_{j,\ell}) > \lambda$ for some pair (j, ℓ) . In words, each Voronoi set \mathcal{V}_i can intersect at most one true cluster \mathbb{B}_s that encloses some Voronoi boundary with a large relative volume.
- Observation 2. For each $s \in [k]$, if $\rho_s(\partial_{j,\ell}) \leq \lambda$ for all pair (j,ℓ) , then $\beta_s^* \in \mathcal{V}_i$ implies that $s \in A_i$. In words, if all Voronoi boundaries enclosed by a true cluster \mathbb{B}_s have small relative volumes, then the center β_s^* cannot itself lie on a Voronoi boundary.

Proof of observations. We prove these observations by contradiction. For Observation [], suppose otherwise that there exists $s \neq s' \in T_i$ for which the statement holds. Part 1 of Theorem [4] ensures that $\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_s^*\| \leq \frac{k^*}{\lambda} + 3r$ and $\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{s'^*}\| \leq \frac{k^*}{\lambda} + 3r$. Using the triangle inequality and the value for λ , we obtain that $\|\boldsymbol{\beta}_s^* - \boldsymbol{\beta}_{s'}^*\| \leq \frac{2k^*}{\lambda} + 6r \leq \Delta_{\max} \frac{8k^*}{c\sqrt{\eta_{\max}}}$, which contradicts the assumption on η_{\min} .

For Observation 2 suppose otherwise that $\beta_s^* \in \mathcal{V}_i$ and $s \notin A_i$ for some $i \in [k]$, which implies that β_s^* lies on a Voronoi boundary and hence $\beta_s^* \notin A_j, \forall j \in [k]$. If $s \in T_j$, then Part 2 of Theorem 4 ensures that $\mathbb{P}_s(\mathcal{V}_j) \leq k\lambda r, \forall j \in [k]$; if $s \notin T_j$, then $\mathbb{P}_s(\mathcal{V}_j) = 0$ by definition of T_j . Summing over $j \in [k]$, we obtain that $1 = \mathbb{P}_s(\mathbb{B}_s) = \sum_{j \in [k]} \mathbb{P}_s(\mathcal{V}_j) \leq k^2 \lambda r < 0.5$, which is a contradiction.

We now construct a partition $\bigcup_{a=0}^{m} S_a = [k]$ of the fitted centers and a partition $\bigcup_{a=1}^{m} S_a^* = [k^*]$ of the true centers that satisfy the conclusion of Theorem 3. These partitions induce an association between the fitted centers in S_a and the true centers in S_a^* , for each $a = 1, \ldots, m$. The construction proceeds in three steps.

a) Step 1 (almost-empty association): First consider the fitted centers indexed by the set

$$\begin{split} S_0 &:= \left\{ i \in [k] : \rho_s(\partial_{j,\ell}) \leq \lambda, \\ \forall (s,j,\ell) \in T_i \times [k] \times [k]; |A_i| = 0 \right\}. \end{split}$$

Part 2(a) of Theorem 4 ensures that for all $i \in S_0$, we have $\mathbb{P}(\mathcal{V}_i) \leq k\lambda r = \frac{ck}{\sqrt{\eta_{\max}}}$ as claimed in Theorem 3.

b) Step 2 (many/one-fit-one association):: We next consider the fitted centers indexed by the set

$$\mathcal{J} := \{ i \in [k] : |A_i| \le 1 \} \setminus S_0.$$

For each $i \in \mathcal{J}$, there are two complementary cases:

ρ_s(∂_{j,ℓ}) ≤ λ for all (s, j, ℓ) ∈ T_i × [k] × [k]; that is, all true clusters that intersect V_i only enclose Voronoi boundaries with a small relative volume. Since i ∉ S₀, by definition of S₀ and J we must have |A_i| = 1; say A_i = {s}. Applying Part 2(b) of Theorem 4, we have

$$\begin{split} \|\beta_{i} - \beta_{s}^{*}\| &= \|\beta_{i} - \mathbf{b}_{i}\| \\ \leq & \frac{k^{*}r}{1 - k^{2}\lambda r} + \frac{k^{*}(k^{2}\lambda r^{2})(1 + (k^{*}k - k^{2})\lambda r)}{(1 - k^{2}\lambda r)^{2}} \\ &+ \frac{(k + 2k^{*})k\lambda r}{1 - k^{2}\lambda r} \Delta_{\max} \\ \stackrel{(i)}{\leq} & 2k^{*}r + 2k^{*}r \left[1 + (k^{*}k - k^{2})\frac{c}{\sqrt{\eta_{\max}}} \right] \\ &+ 2k(k + 2k^{*})\frac{c}{\sqrt{\eta_{\max}}} \Delta_{\max} \\ \stackrel{(ii)}{\leq} & \Delta_{\max}\frac{(4k^{*} + 2ck^{2})(1 + \frac{2k^{*}}{k})}{\sqrt{\eta_{\max}}}, \end{split}$$
(12)

where in step (i) we plug in $\lambda = \frac{c}{\sqrt{\Delta_{\max}r}}$ and use the assumption that $k^2\lambda r < \frac{1}{2}$; in step (ii), we use the assumption that $\eta_{\max} \geq 4c^2k^4$ with c > 3 to further simplify the bound.

ρ_s(∂_{j,ℓ}) > λ for some (s, j, ℓ) ∈ T_i × [k] × [k]; that is, there exists some ground truth cluster B_s that encloses a Voronoi boundary with a large relative volume. Applying Part 1 of Theorem 4 and plugging the value of λ, we obtain that ||β_i − β^{*}_s|| ≤ k^{*}/_λ + 3r ≤ Δ_{max} 4k^{*}/_{c√ηmax}.

In both cases, we have $\|\beta_i - \beta_s^*\| \leq \Delta_{\max} \frac{(4k^* + 2ck^2)(1 + \frac{2k^*}{k})}{\sqrt{\eta_{\max}}}$ as claimed in Theorem 3. For each distinct index $s \in [k^*]$ that appears in the above arguments, set $S_a^* = \{s\}$ and let the corresponding S_a index those β_i 's for which either of the two cases holds. The sets $\{S_a\}$ constructed in this way are disjoint. Indeed, for each i the above two cases are exclusive, where in the first case A_i contains s and only s, and in the second case above the index s is unique by Observation 1.

c) Step 3 (one-fit-many association):: We are left with the fitted centers indexed by the set

$$\mathcal{K} := \{ i \in [k] : |A_i| \ge 2 \} = [k] \setminus (S_0 \cup \mathcal{J}).$$

Similarly to before, for each $i \in \mathcal{K}$, there are two complementary cases:

- $\rho_s(\partial_{j,\ell}) \leq \lambda$ for all $(s, j, \ell) \in T_i \times [k] \times [k]$. Applying Part 2(b) of Theorem 4 and following the same steps as in equation (12), we obtain that $\|\beta_i b_i\| \leq \Delta_{\max} \frac{(4k^* + 2ck^2)(1 + \frac{2k^*}{k})}{\sqrt{\eta_{\max}}}$. In this case, we let $S_a = \{i\}$ and $S_a^* = A_i$. Note that $|S_a^*| = |A_i| \geq 2$ by definition of \mathcal{K} .
- ρ_s(∂_{j,ℓ}) > λ for some (s, j, ℓ) ∈ T_i × [k] × [k]. In this case, applying Part 1 of Theorem 4 would show that β_i is close to β^{*}_s. In fact, we can establish a stronger result showing that β_i is close to the mean of all the true centers contained in its Voronoi set, regardless of whether we include or exclude β^{*}_s. This is the content of the following lemma, which is proved in Section A

Lemma 4 (Proximity to mean of true centers). Under the assumption of Theorem [3] let β be a local minimum

of G. The following is true for each $i \in [k]$. If $\rho_s(\partial_{j,\ell}) > \lambda = \frac{c}{\sqrt{r\Delta_{\max}}}$ for some $(s, j, \ell) \in T_i \times [k] \times [k]$ and $|A_i \setminus \{s\}| \ge 1$, then we have the bounds

$$\begin{aligned} \|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}^{-}\| &\leq \Delta_{\max} \frac{2ck^{2}\left(1 + \frac{2k^{*}}{k}\right) + 4k^{*}\left(3 + \frac{2k^{*}}{k}\right)}{\sqrt{\eta_{\max}}},\\ \|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}^{+}\| &\leq \Delta_{\max} \frac{2ck^{2}\left(1 + \frac{2k^{*}}{k}\right) + 4k^{*}\left(3 + \frac{2k^{*}}{k}\right)}{\sqrt{\eta_{\max}}},\end{aligned}$$

where $\mathbf{b}_i^- := \frac{1}{|A_i \setminus \{s\}|} \sum_{s' \in A_i \setminus \{s\}} \boldsymbol{\beta}_{s'}^*$ and $\mathbf{b}_i^+ := \frac{1}{|A_i \cup \{s\}|} \sum_{s' \in A_i \cup \{s\}} \boldsymbol{\beta}_{s'}^*$; moreover, we have $|A_i \setminus \{s\}| \ge 2$.

In this case, we set $S_a = \{i\}$. Also set $S_a^* = A_i \setminus \{s\}$ if the index s has appeared in the sets $\{S_{a'}^*\}$ constructed previously in Step 2 or in this step, and set $S_a^* = A_i \cup \{s\}$ otherwise. Note that $|S_a^*| \ge 2$ by Lemma 4. As shall become clear momentarily, the flexibility allowed by Lemma 4 is important for ensuring that $\{S_a^*\}$ indeed partitions $[k^*]$.

In both cases above, we have the bound $\|\beta_i - \frac{1}{|S_a^*|} \sum_{s' \in S_a^*} \beta_{s'}^*\| \leq \Delta_{\max} \frac{2ck^2(1+\frac{2k^*}{k})+4k^*(3+\frac{2k^*}{k})}{\sqrt{\eta_{\max}}}$ as claimed in Theorem 3 It is clear that the sets $\{S_a^*\}$ constructed in this step are disjoint from each other and from those constructed in Step 2, because the sets $\{A_i\}$ are disjoint as each true center can be in the interior of only one Voronoi set.

d) Summary: The above procedure constructs a collection of sets $\{S_a\}_{a=0}^m$ and $\{S_a^*\}_{a=1}^m$, which index the fitted and true centers, respectively, and satisfy the bounds in Theorem 3. The sets $\{S_a\}$ indeed form a partition of [k], as we have $\bigcup_{a=0}^m S_a =$ $S_0 \cup \mathcal{J} \cup \mathcal{K} = [k]$ by definition, and $S_a \cap S_b = \emptyset, \forall a \neq b$ by construction and the fact that the three sets $S_0, \mathcal{J}, \mathcal{K}$ are disjoint. For the sets $\{S_a^*\}$, we have argued in the construction above that they are disjoint. On the other hand, each true center β_s^* must belong to at least one Voronoi set \mathcal{V}_i , in which case we have $s \in T_i$. Consider two complementary cases: (i) If $\exists (j, \ell) : \rho_s(\partial_{j, \ell}) > \lambda$, then s must be covered in the second case of either Step 2 or Step 3. (ii) If $\forall (j, \ell) : \rho_s(\partial_{i,\ell}) \leq \lambda$, then Observation 2 ensures that $s \in A_i \subseteq T_i$. In this case, if all other $s' \in T_i$ satisfies $\forall (j, \ell) : \rho_{s'}(\partial_{j,\ell}) \leq \lambda$ as well, then s is covered in the first case of either Step 2 or Step 3. Otherwise, if there exists another $s' \in T_i$ satisfying $\exists (j, \ell) : \rho_{s'}(\partial_{j,\ell}) > \lambda$, then s is covered in the second case of Step 3 (with the role of s and s' exchanged therein) as $s \in A_i \subseteq A_i \setminus \{s'\}$. We conclude that the collection of sets $\{S_a^*\}$ covers all $s \in [k]$ and hence is indeed a partition of [k]. This completes the proof of Theorem 3

C. Proof of Theorem 4

In this section, we prove Theorem $[\underline{4}]$ which shows that a local minimum β satisfies a family of bounds that imply our main Theorem $[\underline{3}]$

a) Proof strategy: To derive structural properties of the local minimum β , we exploit the fact that t = 0 is a local minimum of the directional objective $H^{\boldsymbol{v}}(t)$ (or a smooth upper bound thereof) for any perturbation direction $\boldsymbol{v} \in \mathbb{R}^{d \times k}$; see



Figure 4. Illustration of the quantities $d_{i,j}$, $\theta_{i,j}^{v}$, $D_{i,j,s}$ and $\rho_s(\partial_{i,j})$. The red dots represent β_i and β_j , and the blue star represents the perturbed solution $\beta_j + t v_j$. Here $d_{i,j}$ is the distance between β_i and the mid point $\frac{\beta_i + \beta_j}{2}$ (represented by a green dot); $D_{i,j,s}$, which is represented by the red line segment, is the distance between $\frac{\beta_i + \beta_j}{2}$ and the ball \mathbb{B}_s when computed within the hyperplane containing the Voronoi boundary $\partial_{i,j}$; $\theta_{i,j}^v$ is the angle between the perturbation direction v_j and the direction of $\beta_i - \beta_j$; $\rho_s(\partial_{i,j})$ is the normalized relative volume of the set $\partial_{i,j} \cap \mathbb{B}_s$, which is represented by the blue line segment inside the ball.

Lemma 3 and the discussion in Section ∇ . The expression (10) of $H^{\boldsymbol{v}}$ involves the set $\Delta_{i \to j}^{\boldsymbol{v}}$, which are points that switch from one Voronoi set from another when β is perturbed to $\beta + tv$. For a general direction v, these sets are quite complicated. Our main idea is to focus on a special class of directions satisfying

$$\|\boldsymbol{v}_i\| = 1, \forall i \in [k]; \quad \boldsymbol{v}_i = \boldsymbol{v}_j \text{ or } \boldsymbol{v}_i = -\boldsymbol{v}_j, \forall i \neq j.$$
 (13)

That is, we perturb the β_i 's along the same or opposite directions. For these choices of v, the Voronoi boundary $\partial_{i,j}(\boldsymbol{\beta} + t\boldsymbol{v})$ behaves in a simple way. In particular, when $v_i = v_j$, the boundary $\partial_{i,j}(\beta + tv)$ translates along the direction of v_i ; when $v_i = -v_j$, the boundary rotates around the mid point $\frac{\beta_i + \beta_j}{2}$. See Figure 3 for an illustration. Using this fact, we can construct simple, tractable upper bounds of H^{v} , from which we can deduce the structural properties of the local minimum β .

b) Key quantities: Our analysis involves several key quantities related to the Voronoi sets of β and their boundaries. In particular, for each pair $i \neq j$ whose associated Voronoi sets are adjacent, i.e., $\mathcal{V}_i(\beta) \sim \mathcal{V}_i(\beta)$, we introduce the following four quantities.

- (a) Denote by $d_{i,j} := \frac{1}{2} \|\beta_i \beta_j\|$ the distance between β_i (or β_j) and the Voronoi boundary $\partial_{i,j} \equiv \partial_{i,j}(\beta)$.
- (c) βj and θ^v (i,j) = ∠(v_j, β_i β_j) the angle of the perturbation direction v_j with respect to β_i β_j.
 (c) Define D_{i,j,s} := dist(^{β_i+β_j}/₂, 𝔅_s ∩ ∂_{i,j}), with the convention that D_{i,j,s} = 1 if 𝔅_s ∩ ∂_{i,j} = Ø. Here D_{i,j,s} is the distance between the mid-point ^{β_i+β_j}/₂ and s-th ground truth cluster 𝔅 where the distance is computed within truth cluster \mathbb{B}_s , where the distance is computed within the hyperplane containing the Voronoi boundary $\partial_{i,i}$.
- (d) Recall the quantity $\rho_s(\partial_{i,j}) := \frac{1}{r^d V_d} \operatorname{ReVol}(\partial_{i,j} \cap \mathbb{B}_s)$ defined in the statement of Theorem 4. Note that $\rho_s(\partial_{i,j})$ is the relative volume of the intersection of the Voronoi boundary $\partial_{i,j}$ and the s-th ground truth cluster \mathbb{B}_s , normalized by the volume of \mathbb{B}_s .
- An illustration of these quantities is given in Figure $\frac{4}{4}$

We are now ready to prove Theorem 4. We begin with the upper bound of H^{v} given in equation (11), restated in an equivalent way below:

$$H^{\boldsymbol{v}}(t) \leq U^{\boldsymbol{v}}(t) + \frac{1}{2k^*} \sum_{(i,j):\mathcal{V}_i \sim \mathcal{V}_j} \sum_{s=1}^{k^*} \Big(W_{i \to j,s}^{\boldsymbol{v}}(t) + W_{j \to i,s}^{\boldsymbol{v}}(t) \Big).$$

$$(14)$$

The above inequality holds with equality at t = 0, since by definition $U^{\boldsymbol{v}}(0) = H^{\boldsymbol{v}}(0)$ and $W^{\boldsymbol{v}}_{i \to j,s}(0) = 0, \forall i, j, s.$ Moreover, when β is a local minimum of G, a quick calculation using Lemma 2 shows that

$$\lim_{t \to 0} \frac{1}{t} \left(U^{\boldsymbol{v}}(t) - U^{\boldsymbol{v}}(0) \right) = 0,$$

$$\lim_{t \to 0} \frac{1}{t^2} \left(U^{\boldsymbol{v}}(t) - U^{\boldsymbol{v}}(0) \right) = 1.$$
 (15)

Under the specific choice of the direction v in equation (13), the function $W_{i \to j}^{\boldsymbol{v}} + W_{j \to i}^{\boldsymbol{v}}$ can be further upper bounded, in a small neighborhood of 0, by a smooth function with nice analytical properties. In particular, when the directions $v_i = v_j$ are the same, such an upper bound $\widetilde{W}_{i,j,s}^{v}$ is given in the following proposition, which is proved in Section VI-D,

Proposition 3 (Upper bound, same direction). Let β be a local minimum of G and v satisfy $\|v_i\| = 1, \forall i \in [k]$. If $v_i = v_i$, then $W_{i \to j,s}^{v} + W_{j \to i,s}^{v}$ is upper bounded in a neighborhood of 0 by some smooth function $W_{i,j,s}^{\boldsymbol{v}}$ satisfying the following properties:

$$1) W_{i,j,s}^{v}(0) = 0;$$

$$2) \frac{d}{dt} W_{i,j,s}^{v}(t) |_{t=0} = 0;$$

$$3) \lim_{t\to 0} \frac{1}{t^2} \widetilde{W}_{i,j,s}^{v}(t) |_{t=0} = -2\cos^2(\theta_{i,j}^{v}) d_{i,j} \cdot \rho_s(\partial_{i,j}).$$

When the directions $v_i = -v_j$ are opposite and $d \ge 2$, an upper bound $\widehat{W}_{i,j,s}^{v}$ is given in the following proposition, which is proved in Section VI-E.

Proposition 4 (Upper bound, opposite direction). Let β be a local minimum of G and v satisfy $||v_i|| = 1, \forall i \in [k]$. If $v_i = -v_j$ and $v_i, v_j \in \mathcal{L}_{i,j,s}$, then $W_{i \rightarrow j,s}^v + W_{j \rightarrow i,s}^v$ is upper bounded in a neighborhood of 0 by some smooth function $W_{i,j,s}^{\boldsymbol{v}}$ satisfying the following properties:

1)
$$\widehat{W}_{i,j,s}^{v}(0) = 0;$$

2) $\frac{d}{dt} \widehat{W}_{i,j,s}^{v}(t) \mid_{t=0} = 0;$
3) $\lim_{t\to 0} \frac{1}{t^2} \widehat{W}_{i,j,s}^{v}(t) \mid_{t=0} = -2 \frac{D_{i,j,s}^2}{d_{i,j}} \sin^2(\theta_{i,j}^v) \cdot \rho_s(\partial_{i,j}).$

Also note that $W_{i \to j,s}^{\boldsymbol{v}} \leq 0, \forall (i, j, s)$ by Remark 1. For each $s \in [k^*]$ and each un-ordered pair $(i, j) \in [k] \times [k]$ satisfying $\mathcal{V}_i \sim \mathcal{V}_j$, combining equation (14) and the last two propositions give the following smooth upper bound $H_{i,j}^{v}$ of H^{v} :

$$\begin{aligned} H^{\boldsymbol{v}}(t) &\leq \widetilde{H}_{i,j}^{\boldsymbol{v}}(t) := U^{\boldsymbol{v}}(t) + \frac{1}{k^*} \widetilde{W}_{i,j,s}^{\boldsymbol{v}}(t) \mathbb{1}_{\{\boldsymbol{v}_i = \boldsymbol{v}_j\}} \\ &+ \frac{1}{k^*} \widehat{W}_{i,j,s}^{\boldsymbol{v}}(t) \mathbb{1}_{\{\boldsymbol{v}_i = -\boldsymbol{v}_j \in \mathcal{L}_{i,j,s}\}} \end{aligned}$$

which is valid in a neighborhood of 0 and satisfies $H_{i,j}^{\boldsymbol{v}}(0) =$ $H^{\boldsymbol{v}}(0)$. Since t = 0 is a local minimum of $H^{\boldsymbol{v}}$, Lemma 3 ensures that

$$\lim_{t \to 0} \frac{1}{t^2} \left[\widetilde{H}_{i,j}^{\boldsymbol{v}}(t) - \widetilde{H}_{i,j}^{\boldsymbol{v}}(0) \right] \ge 0.$$
(16)

Moreover, by combining equation (15), Proposition 3 and Proposition 4, we obtain that

$$\lim_{t \to 0} \frac{1}{t^2} \left[\widetilde{H}_{i,j}^{\boldsymbol{v}}(t) - \widetilde{H}_{i,j}^{\boldsymbol{v}}(0) \right]$$

=1 - $\frac{2}{k^*} \cos^2(\theta_{i,j}^{\boldsymbol{v}}) d_{i,j} \cdot \rho_s(\partial_{i,j}) \mathbb{1}_{\{\boldsymbol{v}_i = \boldsymbol{v}_j\}}$ (17)
- $\frac{2}{k^*} \frac{D_{i,j,s}^2}{d_{i,j}} \sin^2(\theta_{i,j}^{\boldsymbol{v}}) \rho_s(\partial_{i,j}) \mathbb{1}_{\{\boldsymbol{v}_i = -\boldsymbol{v}_j \in \mathcal{L}_{i,j,s}\}}.$

Since equation (17) holds for any choice of v satisfying the condition (13), we may choose v judiciously to simplify the right hand side of (17). By doing so we can show that for each $s \in [k^*]$ and each pair $i \neq j \in [k]$ satisfying $\mathcal{V}_i \sim \mathcal{V}_j$, there hold the inequalities

$$d_{i,j} \cdot \rho_s(\partial_{i,j}) \le \frac{k^*}{2}$$
 and (18)

$$\frac{D_{i,j,s}^2}{d_{i,j}} \cdot \rho_s(\partial_{i,j}) \le \frac{k^*}{2},\tag{19}$$

where the second inequality is valid when $d \ge 2$. To prove the inequality (18), suppose otherwise that $d_{i,j} \cdot \rho_s(\partial_{i,j}) > \frac{k^*}{2}$ for some (i, j, s). We can choose the directions $v_i = v_j = \frac{\beta_i - \beta_j}{\|\beta_i - \beta_j\|}$, which satisfies $\theta_{i,j}^v = 0$. Combining with equation (17) gives

$$\lim_{t\to 0}\frac{1}{t^2}\left(\widetilde{H}^{\boldsymbol{v}}_{i,j}(t)-\widetilde{H}^{\boldsymbol{v}}(0)\right)<0,$$

which contradicts the inequality (16). Similarly, to prove the inequality (19), suppose otherwise that $\frac{D_{i,j,s}^2}{d_{i,j}} \cdot \rho_s(\partial_{i,j}) > \frac{k^*}{2}$ for some (i, j, s) when $d \ge 2$. We can choose v_i and v_j to be two unit vectors in the two-dimensional plane $\mathcal{L}_{i,j,s}$ such that $v_i = -v_j$ and $v_i \perp (\beta_j - \beta_i)$, which satisfies $\theta_{i,j}^v = \frac{\pi}{2}$. Combining with equation (17) gives

$$\lim_{t \to 0} \frac{1}{t^2} \left(\widetilde{H}_{i,j}^{\boldsymbol{v}}(t) - \widetilde{H}^{\boldsymbol{v}}(0) \right) < 0$$

which again contradicts the inequality (16).

In the remaining of the proof, fix an index $i \in [k]$ and a number $\lambda > 0$. We shall use equations (18) and (19) to derive the structural properties of β_i and its Voronoi set \mathcal{V}_i . To this end, we consider two complementary cases that correspond to Part 1 and Part 2 of Theorem 4. Recall that $T_i := \{s \in [k^*] : \mathcal{V}_i \cap \mathbb{B}_s \neq \emptyset\}$.

Case 1: there exists some $(s, j, \ell) \in T_i \times [k] \times [k]$ such that $\rho_s(\partial_{j,\ell}) > \lambda$.: In this case, the Voronoi set \mathcal{V}_i intersects a true cluster \mathbb{B}_s that encloses some Voronoi boundary $\partial_{j,\ell}$ with a large relative volume. Note that this case corresponds to Part 1 of Theorem [4].

Under the case condition, the inequality (18) implies that $d_{j,\ell} \leq \frac{k^*}{2\rho_s(\partial_{j,\ell})} \leq \frac{k^*}{2\lambda}$; equivalently,

$$\left\|\frac{\boldsymbol{\beta}_j + \boldsymbol{\beta}_\ell}{2} - \boldsymbol{\beta}_j\right\| = \left\|\frac{\boldsymbol{\beta}_j + \boldsymbol{\beta}_\ell}{2} - \boldsymbol{\beta}_\ell\right\| \le \frac{k^*}{2\lambda}.$$

We consider the one-dimensional and high-dimensional cases separately. When d = 1, the case condition $\rho_s(\partial_{j,\ell}) > \lambda$ further implies that $\frac{\beta_j + \beta_\ell}{2} \in \mathbb{B}_s$ and hence $|\frac{\beta_j + \beta_\ell}{2} - \beta_s^*| \le r$. It follows that $|\beta_j - \beta_s^*| \le |\beta_j - \frac{\beta_j + \beta_\ell}{2}| + |\frac{\beta_j + \beta_\ell}{2} - \beta_s^*| \le \frac{k^*}{2\lambda} + r$. When $d \ge 2$, we have $D_{j,\ell,s} \le \frac{k^*}{2\lambda}$ by multiplying

the inequalities (18) and (19). Let $z \in \mathbb{B}_s \cap \partial_{j,\ell}$ be the point that attains $\left\|\frac{\beta_j + \beta_\ell}{2} - z\right\| = \text{dist}\left(\frac{\beta_j + \beta_\ell}{2}, \mathbb{B}_s \cap \partial_{j,\ell}\right) = D_{j,\ell,s}$. It follows that

$$\begin{split} \|\boldsymbol{\beta}_{j} - \boldsymbol{\beta}_{s}^{*}\| &\leq \left\|\boldsymbol{\beta}_{j} - \frac{\boldsymbol{\beta}_{j} + \boldsymbol{\beta}_{\ell}}{2}\right\| + \left\|\frac{\boldsymbol{\beta}_{j} + \boldsymbol{\beta}_{\ell}}{2} - \boldsymbol{z}\right\| + \|\boldsymbol{z} - \boldsymbol{\beta}_{s}^{*}\| \\ &\leq \frac{k^{*}}{2\lambda} + \frac{k^{*}}{2\lambda} + r = \frac{k^{*}}{\lambda} + r. \end{split}$$

In either case of d, we have the bound $\|\beta_j - \beta_s^*\| \le \frac{k^*}{\lambda} + r$. Since $s \in T_i$, there exist a point $x \in \mathbb{B}_s \cap \mathcal{V}_i$. It follows that

$$\begin{split} \|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{s}^{*}\| &\leq \|\boldsymbol{\beta}_{i} - \boldsymbol{x}\| + \|\boldsymbol{x} - \boldsymbol{\beta}_{s}^{*}\| \\ &\stackrel{(i)}{\leq} \|\boldsymbol{\beta}_{j} - \boldsymbol{x}\| + \|\boldsymbol{x} - \boldsymbol{\beta}_{s}^{*}\| \\ &\leq (\|\boldsymbol{\beta}_{j} - \boldsymbol{\beta}_{s}^{*}\| + \|\boldsymbol{x} - \boldsymbol{\beta}_{s}^{*}\|) + \|\boldsymbol{x} - \boldsymbol{\beta}_{s}^{*}\| \stackrel{(ii)}{\leq} \frac{k^{*}}{\lambda} + 3r, \end{split}$$

where step (i) follows from $x \in \mathcal{V}_i$, and step (ii) follows from $x \in \mathbb{B}_s$ and the bound on $\|\beta_j - \beta_s^*\|$ proved above. We have established Part 1 of Theorem 4.

Case 2: for all $(s, j, \ell) \in T_i \times [k] \times [k]$ there holds $\rho_s(\partial_{j,\ell}) \leq \lambda$.: In this case, for all true clusters $\{\mathbb{B}_s\}$ that intersect the Voronoi set \mathcal{V}_i , all the Voronoi boundaries enclosed by \mathbb{B}_s have a small relative volume.

Let us partition the set $T_i := \{s \in [k^*] : \mathcal{V}_i \cap \mathbb{B}_s \neq \emptyset\}$ into two subsets defined as follows:

$$A_i := \{ s \in T_i : \beta_s^* \in \operatorname{int}(\mathcal{V}_i) \} \text{ and} \\ B_i := T_i \setminus A_i = \{ s \in T_i : \beta_s^* \notin \operatorname{int}(\mathcal{V}_i) \}.$$

Here A_i indexes the ground truth clusters whose centers are in the interior of the Voronoi set \mathcal{V}_i ; B_i indexes the ground truth clusters that intersect \mathcal{V}_i but their centers are outside its interior (i.e., the center either lies on a Voronoi boundary or in some other Voronoi set \mathcal{V}_j). Also recall the quantities $m_{i,s} \equiv m_{i,s}(\beta)$ and $c_{i,s} \equiv c_{i,s}(\beta)$ introduced after Lemma in particular, $m_{i,s}$ is the probability mass of $\mathcal{V}_i \cap \mathbb{B}_s$ with respect to \mathbb{P}_s , and $c_{i,s}$ is the corresponding center of mass.

Note the Voronoi set \mathcal{V}_i is a polyhedron with at most k facets. For each $s \in B_i$, if $\rho_s(\partial_{j,\ell}) \leq \lambda$ for all $\forall (j,\ell)$ (and in particular, for j = i), then all facets of \mathcal{V}_i intersect the ball \mathbb{B}_s with a small relative volume. Moreover, β_s^* is not in $int(\mathcal{V}_i)$. With these two facts, an elementary geometric argument (formally given in Lemma \mathfrak{V}) shows that the intersection $\mathcal{V}_i \cap \mathbb{B}_s$ must have a small mass; that is,

$$m_{i,s} = \mathbb{P}_s(\mathcal{V}_i) \le k\lambda r, \quad \forall s \in B_i.$$
 (20)

On the other hand, for each $s \in A_i$, we must have $\beta_s^* \notin int(\mathcal{V}_j), \forall j \neq i$, since the interiors of Voronoi sets are disjoint. Repeating the same argument above shows that $\mathbb{P}_s(\mathcal{V}_j) \leq k\lambda, \forall j \neq i$, whence

$$m_{i,s} = \mathbb{P}_s(\mathcal{V}_i) = 1 - \sum_{j:j \neq i} \mathbb{P}_s(\mathcal{V}_j) \ge 1 - k^2 \lambda r, \quad \forall s \in A_i.$$
(21)

The inequalities (20) and (21) establish the first two bounds in Part 2 of Theorem 4

We next turn to Part 2(a) of Theorem 4 which concerns the case with $A_i = \emptyset$. This means that $T_i = B_i$. We therefore have

$$\mathbb{P}(\mathcal{V}_i) = \frac{1}{k^*} \sum_{s \in [k^*]} m_{i,s} \le k \lambda r,$$

where the last step holds due to equation (20) and the fact that $m_{i,s} = 0, \forall s \notin T_i$.

Finally, we consider Part 2(b) of Theorem 4 which concerns the case with $A_i \neq \emptyset$ and hence $\mathbb{P}(\mathcal{V}_i) > 0$. Since β is a local minimum, Lemma 2 and the discussion thereafter ensure that

$$\beta_{i} = \frac{\sum_{s=1}^{k^{*}} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s=1}^{k^{*}} m_{i,s}} = \frac{\sum_{s \in T_{i}} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in T_{i}} m_{i,s}}$$

Recalling the definition $b_i := \frac{1}{|A_i|} \sum_{s \in A_i} \beta_s^*$, we can decompose the quantity $(\beta_i - b_i)$ of interest as follows:

$$\beta_{i} - \boldsymbol{b}_{i} = \frac{\sum_{s \in A_{i}} m_{i,s} \boldsymbol{c}_{i,s} + \sum_{s \in B_{i}} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in A_{i}} m_{i,s} + \sum_{s \in B_{i}} m_{i,s}} - \frac{1}{|A_{i}|} \sum_{s \in A_{i}} \beta_{s}^{*}$$

$$= \underbrace{\frac{\sum_{s \in A_{i}} m_{i,s} \boldsymbol{c}_{i,s} + \sum_{s \in B_{i}} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in A_{i}} m_{i,s} + \sum_{s \in B_{i}} m_{i,s}}}_{\boldsymbol{\mu}} - \underbrace{\frac{\sum_{s \in A_{i}} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in A_{i}} m_{i,s}}}_{\sum_{s \in A_{i}} m_{i,s}} - \frac{1}{|A_{i}|} \sum_{s \in A_{i}} \beta_{s}^{*}}.$$
(22)

The following two lemmas, proved in Section VI-F, control the norms of the vectors μ and ν .

Lemma 5. We have $\|\mu\| \leq \frac{k^*r}{1-k^2\lambda r} + \frac{k^*\cdot(k\lambda r)\cdot(k^2\lambda r^2)}{(1-k^2\lambda r)^2} + \frac{2k^*k\lambda r}{1-k^2\lambda r}\Delta_{\max}$.

Lemma 6. We have $\|\boldsymbol{\nu}\| \leq \frac{k^2 \lambda r^2}{1-k^2 \lambda r} + \frac{k^2 \lambda r}{1-k^2 \lambda r} \Delta_{\max}$.

Applying these two lemmas to bound the right hand side of equation (22), we obtain that

$$\begin{split} \|\boldsymbol{\beta}_i - \boldsymbol{b}_i\| \leq & \frac{k^*r}{1 - k^2\lambda r} + \frac{k^2\lambda r^2(1 + (k^*k - k^2)\lambda r)}{(1 - k^2\lambda r)^2} \\ & + \frac{(k + 2k^*)k\lambda r}{1 - k^2\lambda r} \Delta_{\max}, \end{split}$$

thereby proving Part 2(b) of Theorem 4

We have completed the proof of Theorem 4.

D. Proof of Proposition 3 (Same Direction)

Under the perturbation direction $v_i = v_j$, the new Voronoi boundary $\partial_{i,j}(\beta + tv)$ is a translation of the original boundary $\partial_{i,j}(\beta)$ by the amount tv; see top panel of Figure 3 When $\partial_{i,j}(\beta)$ intersects \mathbb{B}_s trivially, with measure 0 with respect to \mathbb{P}_s , setting $\widetilde{W}_{i,j,s}^v \equiv 0$ satisfies the conclusions of the proposition as $\rho_s(\partial_{i,j}) = 0$ in this case. Thus we only need to consider the case where $\partial_{i,j}(\beta)$ intersects \mathbb{B}_s non-trivially. We assume WLOG that $\mathbb{P}_s(\Delta_{i\to j}^v) > 0$ and upper bound $W_{i\to j,s}^v + W_{j\to i,s}^v$ by $W_{i\to j,s}^v$ (as $W_{j\to i,s}^v$ is non-positive). It remains to upper bound $W_{i\to j,s}^v$.

Recall expression for $W_{i \to j,s}^{\boldsymbol{v}}$:

$$egin{aligned} W^{m{v}}_{m{v}
ightarrow j,s}(t) \ &:= \int_{\Delta^{m{v}}_{i
ightarrow j}(t)} (\|m{x}-m{eta}_j-tm{v}_j\|^2-\|m{x}-m{eta}_i-tm{v}_i\|^2)f_s(m{x})\mathrm{d}m{x} \ &= \int_{\Delta^{m{v}}_{i
ightarrow j}(t)} \left[2\langlem{x},m{eta}_i+tm{v}_i-m{eta}_j-tm{v}_j
ight
angle \ &+ (\|m{eta}_j+tm{v}_j\|^2-\|m{eta}_i+tm{v}_i\|^2)
ight]f_s(m{x})\mathrm{d}m{x}. \end{aligned}$$

Since the integrand above only involves the Euclidean norm, we are free to choose any coordinate system. In particular, we choose the origin to be $\frac{1}{2}(\beta_i + \beta_j)$, the principal axis to be the direction of $\beta_i - \beta_j$, and the secondary axis to be the direction orthogonal to $\beta_i - \beta_j$ and in span $\{\beta_i - \beta_j, v_i\}$. Under this coordinate system, we have

$$W_{i \to j,s}^{\boldsymbol{v}}(t) = 2 \int_{\Delta_{i \to j}^{\boldsymbol{v}}(t)} (\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| x_1 - t \langle \boldsymbol{\beta}_i - \boldsymbol{\beta}_j, \boldsymbol{v}_i \rangle) f_s(\boldsymbol{x}) d\boldsymbol{x}$$

$$= 2 \int_{x_2, \dots, x_d : \boldsymbol{x} \in \Delta_{i \to j}^{\boldsymbol{v}}(t)} \int_{x_1 = 0}^{t \cos(\theta)} \left[\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| x_1 - t \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| \cos(\theta) \right] f_s(\boldsymbol{x}) dx_1 dx_2 \dots dx_d,$$

where we introduce the shorthand $\theta \equiv \theta_{i,j}^v := \angle (\beta_i - \beta_j, v_j) = \angle (\beta_i - \beta_j, v_i)$ and, slightly abusing notation, still use f_s to denote the density function under the new coordinate system. Note that the region

$$S(z,t) := \{(x_2,\ldots,x_d) : \boldsymbol{x} \in \Delta_{i \to j}^{\boldsymbol{v}}(t), x_1 = z\}$$

is a vertical slice under the current coordinate system; in particular, S(z,t) is the intersection of the set $\Delta_{i \to j,s}^{v}(t)$ and the hyperplane that is parallel to $\partial_{i,j}$ and at a distance z from $\partial_{i,j}$. Defining the integral

$$\rho_{i \to j,s}^{\boldsymbol{v}}(z,t) := \int_{S(z,t)} f_s(z, x_2, \dots, x_d) \mathrm{d} x_2 \dots \mathrm{d} x_d,$$

we can write

$$W_{i \to j,s}^{\boldsymbol{v}}(t) = 2 \int_{x_1=0}^{t \cos(\theta)} \left[\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| x_1 - t \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| \cos(\theta) \right] \rho_{i \to j,s}^{\boldsymbol{v}}(x_1, t) \mathrm{d}x_1.$$
(23)

When t is small, we have the sandwich bound $m(t) \leq \rho_{i \to j,s}^{v}(x_1, t) \leq M(t)$, where

 $m(t) := \min_{x_1 \in [0, t \cos(\theta)]} \rho_{i \to j, s}^{\boldsymbol{v}}(x_1, t),$

and

$$M(t) := \max_{x_1 \in [0, t \cos(\theta)]} \rho_{i \to j, s}^{\boldsymbol{v}}(x_1, t)$$

Here m(t) and M(t) are well-defined as they are the max/min of the bounded function $\rho_{i \to j,s}^{v}$ over the compact interval $[0, t \cos(\theta)]$. Moreover, m(t) and M(t) satisfy

$$\lim_{t \to 0} m(t) = \lim_{t \to 0} M(t) = \frac{\operatorname{ReVol}(\partial_{i,j})}{\operatorname{Vol}(\mathbb{B}_s(r))} = \rho_s(\partial_{i,j}).$$

Bounding the two terms in the bracket in equation (23) separately, we obtain that

$$2\int_{x_1=0}^{t\cos(\theta)} \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| x_1 \rho_{i \to j,s}^{\boldsymbol{v}}(x_1, t) \mathrm{d}x_1$$

$$\in \left[2d_{i,j}\cos^2(\theta) \cdot m(t)t^2, 2d_{i,j}\cos^2(\theta) \cdot M(t)t^2 \right]$$

and

$$2t \int_{x_1=0}^{t\cos(\theta)} \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| \cos(\theta) \rho_{i \to j,s}^{\boldsymbol{v}}(x_1, t) \mathrm{d}x_1$$

$$\in \left[4d_{i,j}\cos^2(\theta) \cdot m(t)t^2, 4d_{i,j}\cos^2(\theta) \cdot M(t)t^2\right]$$



Figure 5. Illustration of the local coordinate system and the upper bound function. The local coordinate system has the origin at $\frac{\beta_i + \beta_j}{2}$, represented by the dark green dot. Its principal axis is in the direction of $\beta_i - \beta_j$, plotted as the x_1 axis, and its secondary axis is plotted as the x_2 axis. In the left panel, the red dots represent β_i and β_j respectively; the blue stars represent $\beta_i + tv_i$ and $\beta_j + tv_j$ respectively, with $v_i = -v_j$. The dark blue arrow indicates the direction of v_j and it has an angle $\theta_{i,j}^v$ with the vector $\beta_i - \beta_j$, the x_1 axis. Correspondingly, the Voronoi boundary $\partial_{i,j}(\beta + tv)$ rotates around the origin with an angle $\psi(t)$. The boundaries $\partial(\beta + tv)$ are plotted using dotted lines. The shaded green region in the left panel corresponds to the set $\Delta_{i \to j}^v(t)$, in which the point becomes closer to $\beta_j + tv_j$ than $\beta_i + tv_i$ after β is moved to $\beta + tv$. In the right panel, we demonstrate the set $\Delta_{i \to j}^v(t)$ using the shaded green region. It is a subset of $\Delta_{i \to j}^v(t)$, enclosed by the hyperplane $\{x : x_1 = 0\}$ and the translated hyperplane $\{x : x_1 = D_{i,j,s} \tan(\psi(t))\}$

whence

$$2(m(t) - 2M(t))d_{i,j}\cos^2(\theta)t^2 \le W^{\boldsymbol{v}}_{i \to j,s}(t)$$
$$\le 2(M(t) - 2m(t))d_{i,j}\cos^2(\theta)t^2.$$

It is then easy to see that $W_{i \to j,s}^{\boldsymbol{v}}(0) = 0$, $\frac{d}{dt}W_{i \to j,s}^{\boldsymbol{v}}(t) \mid_{t=0} = 0$ and

$$\lim_{t \to 0} \frac{W_{i \to j,s}^{s}(t)}{t^{2}} = -2d_{i,j}\cos^{2}(\theta)\rho_{s}(\partial_{i,j}).$$

In summary, setting $\widetilde{W}_{i,j,s}^{\boldsymbol{v}} = W_{i \to j,s}^{\boldsymbol{v}}$, we have established that $W_{i \to j,s}^{\boldsymbol{v}} + W_{j \to i,s}^{\boldsymbol{v}} \leq \widetilde{W}_{i,j,s}^{\boldsymbol{v}}$ and that $\widetilde{W}_{i,j,s}^{\boldsymbol{v}}$ satisfies the desired analytical properties in Proposition 3.

E. Proof of Proposition 4 (Opposite Direction)

Under the perturbation direction $v_i = -v_j$, the Voronoi boundary $\partial_{i,j}(\beta)$ rotates around the mid point $\frac{\beta_i + \beta_j}{2}$; see right panel of Figure 3. When $\partial_{i,j}(\beta)$ intersects \mathbb{B}_s trivially, with measure 0 with respect to \mathbb{P}_s , setting $\widehat{W}_{i,j,s} \equiv 0$ satisfies the conclusion of the proposition, as $\rho_s(\partial_{i,j}) = 0$ in this case. When $\frac{1}{2}(\beta_i + \beta_j) \in \mathbb{B}_s$, we can also set $\widehat{W}_{i \to j,s}^v \equiv 0$, as $D_{i,j,s} = 0$ in this case. Thus in the rest of the proof, we assume WLOG that $\frac{1}{2}(\beta_i + \beta_j) \notin \mathbb{B}_s$ and $\mathbb{P}_s(\Delta_{i \to j,s}^v(t)) > 0$. In this case we have $\rho_s(\partial_{i,j}) > 0$ and $D_{i,j,s} > 0$. We upper bound $W_{i \to j,s}^v + W_{j \to i,s}^v$ by $W_{i \to j,s}^v$ (as $W_{j \to i,s}^v \leq 0$). It remains to find an upper bound of $W_{i \to j,s}^v$ that satisfies the desired analytical properties.

Similarly to Section VI-D, a convenient coordinate system is used. In particular, we choose the origin to be $\frac{1}{2}(\beta_i + \beta_j)$, the principal axis to be the direction of $\beta_i - \beta_j$, and the secondary axis to be the direction that is orthogonal to $\beta_i - \beta_j$ and in the plane $\mathcal{L}_{i,j,s}$; see the left panel of Figure 5. Under this coordinate system, we have the representation

$$\beta_i = (d_{i,j}, 0, \dots, 0), \quad \beta_j = (-d_{i,j}, 0, \dots, 0),$$

$$\beta_s^* = (b_1^*, b_2^*, 0, \dots, 0)$$
(24)

for some $b_1^*, b_2^* \in \mathbb{R}$. The orientation of the secondary axis can be chosen to satisfy $b_2^* < 0$. We note that the boundary $\partial_{i,j}(\beta)$ rotates around $\frac{1}{2}(\beta_i + \beta_j)$ by an angle $\psi(t)$ that satisfies

$$\tan(\psi(t)) = \frac{t\sin(\theta)}{d_{i,j} - t\cos(\theta)},$$
(25)

where $\theta \equiv \theta_{i,j}^{v} \in [0, \pi/2]$ is the (unsigned) angle between v_j and $\beta_i - \beta_j$. Moreover, since we assume $\mathbb{P}_s(\Delta_{i \to j}^{v}(t)) > 0$, the directions v_i and v_j have the following coordinate representation:

$$oldsymbol{v}_j = ig(\cos(heta), -\sin(heta), 0, \dots, 0ig) = -oldsymbol{v}_i$$

We now proceed to upper bound the function $W_{i \rightarrow j,s}^{v}$. Define the polyhedron set

$$\widetilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t) := \left\{ \boldsymbol{x} \in \Delta_{i \to j}^{\boldsymbol{v}}(t) : x_1 \le D_{i,j,s} \tan(\psi(t)) \right\}.$$
 (26)

The set $\Delta_{i\to j}^{\boldsymbol{v}}(t)$ is sandwiched between the two hyperplanes $x_1 = 0$ and $x_1 = D_{i,j,s} \tan(\psi(t))$; see the right panel of Figure 5 for an illustration. With the above notations, we can upper bound $W_{i\to j,s}^{\boldsymbol{v}}$ as follows:

$$W_{i \to j,s}^{\boldsymbol{v}} \stackrel{(i)}{\leq} \int_{\tilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t)} \left[2 \langle \boldsymbol{x}, \boldsymbol{\beta}_{i} + t\boldsymbol{v}_{i} - \boldsymbol{\beta}_{j} - t\boldsymbol{v}_{j} \rangle \right. \\ \left. + \left(\| \boldsymbol{\beta}_{j} + t\boldsymbol{v}_{j} \|^{2} - \| \boldsymbol{\beta}_{i} + t\boldsymbol{v}_{i} \|^{2} \right) \right] f_{s}(\boldsymbol{x}) d\boldsymbol{x} \\ \stackrel{(ii)}{=} \int_{\tilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t)} 2 \langle \boldsymbol{x}, \boldsymbol{\beta}_{i} + t\boldsymbol{v}_{i} - \boldsymbol{\beta}_{j} - t\boldsymbol{v}_{j} \rangle f_{s}(\boldsymbol{x}) d\boldsymbol{x} \\ \left. = \int_{\tilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t)} \left[2 \| \boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{j} \| x_{1} - 4t \cos(\theta) x_{1} \right. \\ \left. + 4t \sin(\theta) x_{2} \right] f_{s}(\boldsymbol{x}) d\boldsymbol{x},$$
 (27)

where step (i) holds because the integrand in the definition of $W_{i \to j,s}^{v}$ is non-positive and thus integrating over a smaller set $\widetilde{\Delta}_{i \to j}^{v}(t) \subseteq \Delta_{i \to j}^{v}(t)$ does not decrease the value of the integral, and step (ii) holds since under the current coordinate system, $\beta_i = -\beta_j$ and $v_i = -v_j$.

To proceed, we let

$$D(t) := \max\left\{x_2 : \boldsymbol{x} \in \widetilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t) \cap \mathbb{B}_s\right\}.$$
 (28)

denote maximum of the second coordinate of the set $\Delta_{i \to j}^{v}(t) \cap \mathbb{B}_{s}$ under the current coordinate system. The following lemma, proved at the end of this section, characterizes the limit property of D(t).

Lemma 7 (Negative second coordinate at the boundary). Suppose that v satisfies $||v_i|| = ||v_j|| = 1$ and $v_i =$ $-v_j \in \mathcal{L}_{i,j,s}$, $\frac{\beta_i + \beta_j}{2} \notin \mathbb{B}_s$ and $\Delta_{i \to j}^{v}(t) \cap \mathbb{B}_s \neq \emptyset$. We have $\lim_{t \to 0} D(t) = -D_{i,j,s}$.

Lemma 7 ensures that $\lim_{t\to 0} D(t) = -D_{i,j,s} < 0$. Consequently, when t is sufficiently small, we have D(t) < 0 by the continuity.

Continuing from the last display equation (27), we obtain our final upper bound $\widehat{W}_{i,j,s}^{v}(t)$:

$$W_{i \to j,s}^{\boldsymbol{v}} \leq \int_{\tilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t)} \left[2 \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\| x_1 - 4t \cos(\theta) x_1 + 4t \sin(\theta) D(t) \right] f_s(\boldsymbol{x}) d\boldsymbol{x} =: \widehat{W}_{i,j,s}^{\boldsymbol{v}}(t).$$

To establish the analytical properties of $\widehat{W}_{i,j,s}^{\boldsymbol{v}}(t)$, we follow a similar argument as in Section VI-D. Define the integral

$$\rho_{i \to j,s}^{\boldsymbol{v}}(z,t) := \int_{\substack{x_2, \dots, x_d:\\ \boldsymbol{x} \in \tilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t), x_1 = z}} f_s(z, x_2, \dots, x_d) \mathrm{d}x_2 \dots \mathrm{d}x_d,$$

and rewrite $\widehat{W}_{i,j,s}^{\boldsymbol{v}}(t)$ compactly as follows:

$$\widehat{W}_{i,j,s}^{\boldsymbol{v}}(t) = \int_{0}^{D_{i,j,s} \tan(\psi(t))} \left[2 \|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{j}\| x_{1} - 4t \cos(\theta) x_{1} + 4t \sin(\theta) D(t) \right] \rho_{i \to j,s}^{\boldsymbol{v}}(x_{1}, t) \mathrm{d}x_{1}.$$

We have the sandwich bound $m(t) \leq \rho_{i \to j,s}^{\boldsymbol{v}}(x_1,t) \leq M(t)$, valid for $x_1 \in [0, D_{i,j,s} \tan(\psi(t))]$, where $m(t) := \min_{x_1 \in [0, D_{i,j,s} \tanh(\psi(t))]} \rho_{i \to j,s}^{\boldsymbol{v}}(x_1,t)$ and $M(t) := \max_{x_1 \in [0, D_{i,j,s} \tan(\psi(t))]} \rho_{i \to j,s}^{\boldsymbol{v}}(x_1,t)$.

Moreover, the functions $m(\cdot)$ and $M(\cdot)$ satisfy

$$\lim_{t \to 0} m(t) = \lim_{t \to 0} M(t) = \frac{\operatorname{ReVol}(\partial_{i,j})}{\operatorname{Vol}(\mathbb{B}_s)} = \rho_s(\partial_{i,j}).$$

With some algebra as well as the non-positivity of D(t), we obtain the following sandwich bound for $\widehat{W}_{i,j,s}^{\boldsymbol{v}}$ over a small neighborhood of 0:

$$\begin{split} \widehat{W}_{i,j,s}^{\boldsymbol{v}}(t) &\geq 2m(t)d_{i,j}D_{i,j,s}^{2}\tan^{2}(\psi(t)) \\ &+ 4tD_{i,j,s}\sin(\theta)D(t)M(t)\tan(\psi(t)) \\ &- 2tD_{i,j,s}^{2}\cos(\theta)M(t)\tan^{2}(\psi(t)); \end{split}$$

and

$$\begin{split} \bar{W}_{i,j,s}^{\boldsymbol{v}}(t) &\leq 2M(t)d_{i,j}D_{i,j,s}^{2}\tan^{2}(\psi(t)) \\ &+ 4tD_{i,j,s}\sin(\theta)D(t)m(t)\tan(\psi(t)) \\ &- 2tD_{i,j,s}^{2}\cos(\theta)m(t)\tan^{2}(\psi(t)). \end{split}$$

Proposition 4 follows immediately from the limit properties of m(t), M(t), D(t) and $tanh(\psi(t))$.

Proof of Lemma We use the same notations and coordinate system as before. Observe that in equation (28), the maximum that defines D(t) must be attained at a point in $\widetilde{\Delta}_{i\to j}^{v}(t) \cap \mathbb{S}_s$, where \mathbb{S}_s is the hypersphere of the ball \mathbb{B}_s ; see the right panel of Figure 5. We claim that the maximum must also be attained by a point in the hyperplane $\mathcal{L}_{i,j,s}$. Indeed, recalling the representation in equation (24), we see that each point $\boldsymbol{x} \in \widetilde{\Delta}_{i\to j}^{v}(t) \cap \mathbb{S}_s$ must satisfy

$$(x_1 - b_1^*)^2 + (x_2 - b_2^*)^2 + \sum_{j \ge 3} x_j^2 = r^2$$
 and

$$x_1 \in \left[0, D_{i,j,s} \tan(\psi(t))\right]$$

From the above equation it is clear that for each fixed $z \in [0, D_{i,j,s} \tan(\psi(t))]$, over the set $\mathbb{S}_s \cap \widetilde{\Delta}_{i \to j}^{v}(t) \cap \{x : x_1 = z\}$, the maximum of x_2 is attained exactly when $x_j = 0$ for all $j \geq 3$, that is, $x \in \mathcal{L}_{i,j,s}$. Combining these observations, we conclude that

$$D(t) = \max \{ x_2 : \boldsymbol{x} \in \widetilde{\Delta}_{i \to j}^{\boldsymbol{v}}(t) \cap \mathbb{S}_s \cap \mathcal{L}_{i,j,s} \}.$$

Note that the set $\Delta_{i\to j}^{\boldsymbol{v}}(t) \cap \mathbb{S}_s \cap \mathcal{L}_{i,j,s}$ is compact, which is represented by the solid blue segment in \mathbb{B}_s in the right panel of Figure 5; as $t \to 0$, this set shrinks continuously to a single point $(0, -D_{i,j,s}, \ldots, 0)$. It thus follows that $\lim_{t\to 0} D(t) = -D_{i,j,s}$. This completes the proof of Lemma 7.

F. Proofs of Lemmas 5 and 6

For convenience we restate the bounds in equations (20) and (21) as follows:

$$\begin{cases} m_{i,s} \le k\lambda r, & \text{if } s \in B_i, \\ m_{i,s} \ge 1 - k^2\lambda r, & \text{if } s \in A_i, \end{cases}$$
(29)

where we recall that $m_{i,s} = \mathbb{P}(\mathcal{V}_i)$ is the probability mass of the set \mathcal{V}_i with respect to the uniform distribution on the ball $\mathbb{B}_s \equiv \mathbb{B}_{\boldsymbol{\beta}_s^*}(r)$, and \boldsymbol{c}_s is the corresponding center of mass. Applying the simple geometric result in Lemma 10, we further obtain the bound

$$\|\boldsymbol{\beta}_{s}^{*} - \boldsymbol{c}_{i,s}\| \leq \frac{r \cdot (1 - m_{i,s})}{m_{i,s}} \leq \begin{cases} \frac{r}{m_{i,s}}, & \text{if } s \in B_{i}, \\ \frac{k^{2} \lambda r^{2}}{m_{i,s}}, & \text{if } s \in A_{i}, \end{cases}$$
(30)

where the last step follows from the fact that $m_{i,s} \leq 1$ and equation (29). We are ready to prove the two lemmas.

Proof of Lemma 5 We have the following decomposition of the vector μ :

$$\boldsymbol{\mu} = \frac{\sum_{s \in A_i} m_{i,s} \boldsymbol{c}_{i,s} + \sum_{s \in B_i} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s}} - \frac{\sum_{s \in A_i} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s}}$$

$$= \frac{\sum_{s \in B_i} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s}}$$

$$= \frac{\sum_{s \in B_i} m_{i,s} \sum_{s \in A_i} m_{i,s} \sum_{s \in A_i} m_{i,s} \boldsymbol{c}_{i,s}}{\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s}}$$

$$= \frac{\sum_{s \in B_i} m_{i,s} (\boldsymbol{c}_{i,s} - \boldsymbol{\beta}_s^*)}{\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s}}$$

$$- \frac{\sum_{s \in B_i} m_{i,s} \sum_{s \in A_i} m_{i,s} \sum_{s \in A_i} m_{i,s}}{(\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s}) \sum_{s \in A_i} m_{i,s}}$$

$$+ \frac{\sum_{s \in B_i} m_{i,s} \boldsymbol{\beta}_s^*}{\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s}}$$

$$- \frac{\sum_{s \in B_i} m_{i,s} \sum_{s \in A_i} m_{i,s}}{(\sum_{s \in A_i} m_{i,s} + \sum_{s \in B_i} m_{i,s})}$$

It follows that

$$\|\boldsymbol{\mu}\| \leq \frac{\sum_{s \in B_{i}} m_{i,s} \|\boldsymbol{c}_{i,s} - \boldsymbol{\beta}_{s}^{*}\|}{\sum_{s \in A_{i}} m_{i,s} + \sum_{s \in B_{i}} m_{i,s}} + \frac{\sum_{s \in B_{i}} m_{i,s} \sum_{s \in A_{i}} m_{i,s} \|\boldsymbol{c}_{i,s} - \boldsymbol{\beta}_{s}^{*}\|}{(\sum_{s \in A_{i}} m_{i,s} + \sum_{s \in B_{i}} m_{i,s}) \sum_{s \in A_{i}} m_{i,s}}$$

$$+ \frac{2\sum_{s\in B_i} m_{i,s}}{\sum_{s\in A_i} m_{i,s} + \sum_{s\in B_i} m_{i,s}} \Delta_{\max}.$$

$$\stackrel{\text{ij}}{\leq} \frac{k^*r}{1 - k^2\lambda r} + \frac{k^* \cdot (k\lambda r) \cdot (k^2\lambda r^2)}{(1 - k^2\lambda r)^2} + \frac{2k^*k\lambda r}{1 - k^2\lambda r} \Delta_{\max},$$

where the last step holds due to the inequalities (29) and (30) as well as the fact that $|A_i|, |B_i| \in [1, k^*]$. This completes the proof of Lemma 5.

Proof of Lemma 6 We have the following decomposition of the vector ν :

$$\begin{split} \mathbf{v} &= \frac{\sum_{s \in A_i} m_{i,s} \mathbf{c}_{i,s}}{\sum_{s \in A_i} m_{i,s}} - \frac{1}{|A_i|} \sum_{s \in A_i} \boldsymbol{\beta}_s^* \\ &= \frac{\sum_{s \in A_i} m_{i,s} (\mathbf{c}_{i,s} - \boldsymbol{\beta}_s^*)}{\sum_{s \in A_i} m_{i,s}} \\ &+ \frac{\sum_{s \in A_i} (m_{i,s} - \frac{1}{|A_i|} \sum_{s' \in A_i} m_{i,s'}) \boldsymbol{\beta}_s^*}{\sum_{s \in A_i} m_{i,s}} \end{split}$$

Using the inequalities (29) and (30) as well as the triangle inequality, we obtain that

$$\|\boldsymbol{\nu}\| \leq rac{k^2 \lambda r^2}{1 - k^2 \lambda r} + rac{k^2 \lambda r}{1 - k^2 \lambda r} \Delta_{\max},$$

thereby proving Lemma 6.

v

VII. CONCLUSION

In this paper, we characterize the structures of all local minima in the k-means problem for general values of k. We show that under an appropriate separation condition of the ground truth clusters, the local minima are always composed of one-fit-many, many-fit-one or almost-empty type associations between the fitted and ground truth centers.

Several future directions are of interests for both theory and applications. An immediate direction is to generalize our results from the population case to the finite sample case, and from balanced spherical GMMs to more general mixture models with imbalanced clusters, general covariance matrices and heavy-tailed distributions. To transfer the population case to the finite sample case, we might be able to adapt a set of general techniques based on localization and uniform concentration arguments; see [44, 54]. Also, while we have focused on the k-means formulation, we expect that similar structural results hold for a much broader class of clustering formulations, including the maximum likelihood formulation of mixture problems. On the computational side, we have discussed the implications of our results for improving clustering algorithms. Rigorously justifying these algorithms (which are largely heuristic so far) in a broad range of models would be interesting.

Finally, it would be of great interest to establish similar structural results for other non-convex optimization problems that arise in machine learning and statistics applications [55], 46, 47], 45], 44].

APPENDIX

EQUIVALENCE TO THE PARTITION-BASED FORMULATION

A common way of formulating the k-means clustering problem is as follows: given a set of observations $x_1, \ldots, x_n \in \mathbb{R}^d$, we find a partition $S = \{S_1, \ldots, S_k\}$ of these observations such that the within-cluster sum of squared distances is minimized:

$$\min_{\boldsymbol{S}} \sum_{j=1}^{\kappa} \sum_{\boldsymbol{x} \in S_j} \|\boldsymbol{x} - \boldsymbol{\mu}_j\|^2,$$
(31)

where $\mu_j = \frac{1}{|S_j|} \sum_{x \in S_j} x$ is the mean of points in cluster *i*. Meanwhile, the formulation (1) used in this paper, restated below,

$$\min_{oldsymbol{eta}} \sum_{i=1}^n \min_{j \in [k]} \|oldsymbol{x}_i - oldsymbol{eta}_j\|^2,$$

is based on optimizing over the centers $\beta = (\beta_1, \dots, \beta_k)$. Note that each solution β induces a partition $\{S_j(\beta)\}$ of the observations via the Voronoi diagram. Also note that the sum of squared distances to a set of points is minimized by the mean of these points. Combining this observations, we obtain that for any β :

$$\min_{S} \sum_{j=1}^{k} \sum_{\boldsymbol{x} \in S_{j}} \|\boldsymbol{x} - \boldsymbol{\mu}_{j}\|^{2} \leq \sum_{j=1}^{k} \sum_{\boldsymbol{x} \in S_{j}(\boldsymbol{\beta})} \|\boldsymbol{x} - \boldsymbol{\beta}_{j}\|^{2} \\
= \sum_{i=1}^{n} \min_{j \in [k]} \|\boldsymbol{x}_{i} - \boldsymbol{\beta}_{j}\|^{2}. \quad (32)$$

On the other hand, for any partition $S = \{S_1, \ldots, S_k\}$ of the data points and its corresponding means (μ_1, \ldots, μ_k) , we have

$$\sum_{j=1}^{k} \sum_{\boldsymbol{x} \in S_{j}} \|\boldsymbol{x} - \boldsymbol{\mu}_{j}\|^{2} \ge \sum_{i=1}^{n} \min_{j \in [k]} \|\boldsymbol{x}_{i} - \boldsymbol{\mu}_{j}\|^{2}$$
$$\ge \min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \min_{j \in [k]} \|\boldsymbol{x}_{i} - \boldsymbol{\beta}_{j}\|^{2}.$$
(33)

Taking the minimum over β of both sides of equation (32), and the minimum over S for equation (33), we conclude that the two formulations (31) and (1) have the same optimal values. Moreover, an optimal solution for one formulation induces an optimal solution for the other. Hence these two formulations are equivalent.

PROOF OF PROPOSITION 1

In this section we prove Proposition \square , which states that under the Stochastic Ball Model with $k^* = k$, the ground truth centers β^* is the global minimum of the k-means objective function G.

Proof. We begin by upper bounding the objective value of the ground truth:

$$egin{aligned} G(oldsymbol{eta}^*) &= rac{1}{k}\sum_{s\in[k]}\int\min_{i\in[k]}\|oldsymbol{x}-oldsymbol{eta}_i^*\|^2f_s(oldsymbol{x})\mathrm{d}oldsymbol{x} \ &\leq rac{1}{k}\sum_{s\in[k]}\int\|oldsymbol{x}-oldsymbol{eta}_s^*\|^2f_s(oldsymbol{x})\mathrm{d}oldsymbol{x} \ &\leq r^2, \end{aligned}$$

where the second inequality follows from the fact that each true cluster \mathbb{B}_s has radius r.

Now let β be a global minimum of G. By optimality of β , we have for each $s \in [k]$:

γ

$$egin{aligned} &\mathcal{F}^2 \geq G(oldsymbol{eta})\ &\stackrel{ ext{(i)}}{\geq} rac{1}{k} \int \min_{i \in [k]} \|oldsymbol{x} - oldsymbol{eta}_i\|^2 f_s(oldsymbol{x}) \mathrm{d}oldsymbol{x}\ &\stackrel{ ext{(ii)}}{\geq} rac{1}{k} \int \min_{i \in [k]} \left(rac{1}{2} \|oldsymbol{eta}_s^* - oldsymbol{eta}_i\|^2 - \|oldsymbol{x} - oldsymbol{eta}_s^*\|^2
ight) f_s(oldsymbol{x}) \mathrm{d}oldsymbol{x}\ &\stackrel{ ext{(iii)}}{\geq} rac{1}{2k} \min_{i \in [k]} \|oldsymbol{eta}_s^* - oldsymbol{eta}_i\|^2 - rac{r^2}{k}. \end{aligned}$$

where step (i) holds by ignoring k-1 clusters, step (ii) holds by the inequality $(a-b)^2 \ge \frac{1}{2}a^2 - b^2$, and step (iii) holds because f_s is a probability density. From the above equation we obtain that

$$\min_{i \in [k]} \|\boldsymbol{\beta}_s^* - \boldsymbol{\beta}_i\| < 2\sqrt{k}r, \qquad \forall s \in [k];$$

that is, each true center β_s^* is $2\sqrt{kr}$ -close to at least one β_i . We further observe that each β_i is $2\sqrt{kr}$ -close to at most one true center β_s^* ; otherwise, by the triangle inequality we woud have $\Delta_{\min} \leq ||\beta_s^* - \beta_{s'}^*|| \leq ||\beta_s^* - \beta_i|| + ||\beta_{s'}^* - \beta_i|| < 4\sqrt{kr}$, contradicting the SNR assumption $\eta_{\min} := \frac{\Delta_{\min}}{r} \geq 6\sqrt{k}$. Since the number of β_i 's is equal to that of β_s^* 's, we deduce that each β_i is $2\sqrt{kr}$ -close to *exactly one* β_s^* . Without loss of generality, we may assume that

$$\|\boldsymbol{\beta}_s^* - \boldsymbol{\beta}_s\| < 2\sqrt{k}r, \qquad \forall s \in [k].$$

When the above inequality and the SNR assumption $\eta_{\min} := \frac{\Delta_{\min}}{r} \ge 6\sqrt{k}$ hold, we have for each pairs $(s, s') \in [k] \times [k]$ with $s \neq s'$ and each $\boldsymbol{x} \in \mathbb{B}_s$:

$$\begin{aligned} \|\boldsymbol{x} - \boldsymbol{\beta}_s\| &\leq \|\boldsymbol{x} - \boldsymbol{\beta}_s^*\| + \|\boldsymbol{\beta}_s^* - \boldsymbol{\beta}_s\| \\ &< r + 2\sqrt{kr} \\ &\leq 6\sqrt{kr} - r - 2\sqrt{kr} \\ &< \|\boldsymbol{\beta}_s^* - \boldsymbol{\beta}_{s'}^*\| - \|\boldsymbol{x} - \boldsymbol{\beta}_s^*\| - \|\boldsymbol{\beta}_{s'} - \boldsymbol{\beta}_{s'}^*\| \\ &\leq \|\boldsymbol{x} - \boldsymbol{\beta}_{s'}\|, \end{aligned}$$

which implies that $\mathbb{B}_s \subseteq \mathcal{V}_s(\beta)$. Applying Lemma 2 to the global minimum β , we obtain that

$$\boldsymbol{\beta}_s = \frac{\int_{\mathcal{V}_s(\boldsymbol{\beta})} \boldsymbol{x} f(\boldsymbol{x}) \mathrm{d} \boldsymbol{x}}{\int_{\mathcal{V}_s(\boldsymbol{\beta})} f(\boldsymbol{x}) \mathrm{d} \boldsymbol{x}} = \frac{\int_{\mathbb{B}_s} \boldsymbol{x} f(\boldsymbol{x}) \mathrm{d} \boldsymbol{x}}{\int_{\mathbb{B}_s} f(\boldsymbol{x}) \mathrm{d} \boldsymbol{x}} = \boldsymbol{\beta}_s^*, \qquad \forall s \in [k].$$

thereby proving that β^* is the only global minimum.

PROOF OF PROPOSITION 2

In this section we prove Proposition 2, which states that under the one-dimensional Stochastic Ball Model in Figure 2 with r < 0.4, the solution $\beta = (\beta_1, \beta_2, \beta_3) = (-2 - \frac{r}{2}, -2 + \frac{r}{2}, 1)$ is a local minimum of the k-means objective function G.

Proof. Observe that $\mathcal{V}_1(\beta) = (-\infty, -2], \quad \mathcal{V}_2(\beta) = [-2, \frac{-1+r/2}{2}], \quad \mathcal{V}_3(\beta) = [\frac{-1+r/2}{2}, \infty], \quad \partial_{1,2}(\beta) = -2 \text{ and} \\ \partial_{2,3}(\beta) = \frac{-1+r/2}{2}. \text{ When } r < 0.4, \text{ it is easy to see that for any } \mathbf{b} = (b_1, b_2, b_3) \in \mathbb{R}^3 \text{ in a small neighborhood of } \beta, \text{ the Voronoi boundary } \partial_{2,3}(\mathbf{b}) \text{ remains strictly between } -2+r \text{ and} -r, \text{ and } \partial_{1,2}(\mathbf{b}) \text{ remains strictly between } -2-r \text{ and } -2+r.$

Therefore, for any such b we can explicitly write down its objective value:

$$G(\mathbf{b}) = \int_{-2-r}^{\frac{b_1+b_2}{2}} (x-b_1)^2 dx + \int_{\frac{b_1+b_2}{2}}^{-2+r} (x-b_2)^2 dx + \int_{-r}^{r} (x-b_3)^2 dx + \int_{2-r}^{2+r} (x-b_3)^2 dx.$$

We compute the derivative and Hessian for G at b:

$$\nabla_{\boldsymbol{b}}G = \begin{bmatrix} -2\int_{-2-r}^{\frac{b_1+b_2}{2}} (x-b_1) dx \\ -2\int_{-2+r}^{-2+r} (x-b_2) dx \\ -2\int_{-r}^{r} (x-b_3) dx - 2\int_{2-r}^{2+r} (x-b_3) dx \end{bmatrix}$$

and

$$\nabla_b^2 G = \left[\begin{array}{ccc} \frac{b_1 - b_2}{2} + 2(\frac{b_1 + b_2 + 2 + r}{2}) & \frac{b_1 - b_2}{2} & 0\\ \frac{b_1 - b_2}{2} & \frac{b_1 - b_2}{2} + 2(-2 + r - \frac{b_1 + b_2}{2}) & 0\\ 0 & 0 & 8r \end{array} \right]$$

Evaluating the above expressions at $b_1 = \beta_1 = 2 = \frac{r}{2}$, $b_2 = \beta_2 = -2 + \frac{r}{2}$ and $b_3 = \beta_3 = 1$, we find that the derivative vanishes and the Hessian is positive definite:

$$\begin{split} \nabla_{\boldsymbol{b}} G \big|_{\boldsymbol{b} = \boldsymbol{\beta}} &= 0, \\ \nabla_{\boldsymbol{b}}^2 G \big|_{\boldsymbol{b} = \boldsymbol{\beta}} &= \begin{bmatrix} 1.5r & -0.5r & 0\\ -0.5r & 1.5r & 0\\ 0 & 0 & 8r \end{bmatrix} \succ 0. \end{split}$$

Therefore, β is indeed a local minimum of G.

PROOFS FOR SECTION ∇

In this section, we prove the technical lemmas stated in Section ∇ .

Proof of Lemma 1

Proof. Our goal is to derive the existence and expression of the derivative of the function

$$H^{\boldsymbol{v}}(t) := G(\boldsymbol{\beta} + t\boldsymbol{v}) = \int_{\boldsymbol{x}} \min_{i \in [k]} \|\boldsymbol{x} - \boldsymbol{\beta}_i - t\boldsymbol{v}_i\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x},$$

at t = 0. We make use of the following measure-theoretic version of the Leibniz integral rule.

Proposition 5 (Leibniz's integral rule). Let *T* be an open subset of \mathbb{R} , and *X* be a measure space. Suppose $g: T \times X \to \mathbb{R}$ satisfies the following conditions: (i) $g(t, \mathbf{x})$ is a Lebesgueintegrable function of \mathbf{x} for each $t \in T$; (ii) for almost all $\mathbf{x} \in X$, the partial derivative $\frac{\partial}{\partial t}g(t, \mathbf{x})$ exists for all $t \in T$; (iii) There is an integrable function $\theta: X \to \mathbb{R}$ such that $|\frac{\partial}{\partial t}g(t, \mathbf{x})| \leq \theta(\mathbf{x})$ for all $t \in T$ and almost every $\mathbf{x} \in X$, then we have

$$\frac{\mathrm{d}}{\mathrm{d}t}\int_X g(t,\boldsymbol{x})\mathrm{d}\boldsymbol{x} = \int_X \frac{\partial}{\partial t}g(t,\boldsymbol{x})\mathrm{d}\boldsymbol{x}.$$

We verify the above three conditions in the proposition for $H^{\boldsymbol{v}}$. Without loss of generality, assume that $\|\boldsymbol{v}_i\| \leq 1, \forall i \in [k]$. Let $\Delta := \min_{i \neq j} \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j\|$, which satisfies $\Delta > 0$ by the assumption that $\{\boldsymbol{\beta}_j\}_{j=1}^k$ are pairwise distinct. For condition (i), we see that the function $g(t, \boldsymbol{x}) := \min_{i \in [k]} \|\boldsymbol{x} - \boldsymbol{\beta}_i - t\boldsymbol{v}_i\|^2 f(\boldsymbol{x})$ is integrable in \boldsymbol{x} for each bounded t, since the

density f has bounded second moment. For condition (ii), note that when $t \in T := [-\frac{\Delta}{4}, \frac{\Delta}{4}]$, the perturbed solution $\beta + tv$ remains pairwise disjoint, hence the Voronoi boundary $\partial(\beta + tv)$ has measure 0. For all $t \in T$ and all $x \notin \partial(\beta + tv)$, the minimizer in the definition of $g(t+\epsilon, x)$ remains fixed when $|\epsilon|$ is sufficiently small, hence the partial derivative $\frac{\partial}{\partial t}g(t, x)$ exists at all $t \in T$ and satisfies

$$\begin{aligned} \boldsymbol{x} \in \mathcal{V}_i(\boldsymbol{\beta} + t\boldsymbol{v}) \\ \Longrightarrow \frac{\partial}{\partial t} g(t, \boldsymbol{x}) &= -2 \langle \boldsymbol{v}_i, \boldsymbol{x} - \boldsymbol{\beta}_i - t\boldsymbol{v}_i \rangle f(\boldsymbol{x}). \end{aligned} \tag{34}$$

Finally for condition (iii), for each $x \in \text{support}(f) = \bigcup_{s \in [k]} \mathbb{B}_s(r)$, we have the bound $|\langle v_i, x - \beta_i - tv_i \rangle| \leq \max_{s \in [k]} ||\beta_s^*|| + r + ||\beta_i|| + \frac{\Delta}{4}$ when $t \in T$, hence $|\frac{\partial}{\partial t}g(t, x)|$ is bounded by an integrable function. Applying the Leibniz's integral rule and equation (34), we obtain that

$$egin{aligned} &rac{\mathrm{d}}{\mathrm{d}t}H^{m{v}}(0) = \int_{m{x}}rac{\partial}{\partial t}g(0,m{x})\mathrm{d}m{x} \ &= -\sum_{i=1}^k\int_{\mathcal{V}_i(m{eta})}2\langlem{v}_i,m{x}-m{eta}_i
angle f(m{x})\mathrm{d}m{x} \end{aligned}$$

as claimed.

Proof of Lemma 2

Proof. In view of the decomposition (10) and Remark 1, we have the following upper bound for H^{v} :

$$H^{\boldsymbol{v}}(t) \leq h^{\boldsymbol{v}}(t) := \sum_{i=1}^{k} \int_{\mathcal{V}_{i}(\boldsymbol{\beta})} \|\boldsymbol{x} - \boldsymbol{\beta}_{i} - t\boldsymbol{v}_{i}\|^{2} f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}, \quad (35)$$

which satisfies $h^{\boldsymbol{v}}(0) = H^{\boldsymbol{v}}(0)$. Since $\boldsymbol{\beta}$ is a local minimum of G, we know that t = 0 is local minimum of $H^{\boldsymbol{v}}$ for all \boldsymbol{v} , hence Lemma \mathfrak{Z} ensures that t = 0 is also a local minimum of $h^{\boldsymbol{v}}$.

Suppose that we have $\beta_1 = \beta_2$ and $\mathcal{V}_1(\beta) = \mathcal{V}_2(\beta)$ has a positive measure with respect to f, and that all other $\beta_j, j \ge 3$ are pairwise distinct and different from β_1 and β_2 . We may partition $\mathcal{V}_1(\beta) = \mathcal{V}_2(\beta)$ into two disjoint sets S_1 and S_2 , each with positive measure. For $i \in \{1, 2\}$ denote by $s_i := \frac{\int_{S_i} xf(x)dx}{\int_{S_i} f(x)dx}$ the center of mass of S_i with respect to f. We can choose the partition in such a way that $s_1 \neq \beta_1$ and $s_2 \neq \beta_2$. Fix a direction $v = (v_1, v_2, 0, \dots, 0)$ with $v_1 = s_1 - \beta_1$ and $v_2 = s_2 - \beta_2$. In this case the upper bound h^v can be written as

$$\begin{split} h^{\boldsymbol{v}}(t) &= \int_{S_1} \|\boldsymbol{x} - \boldsymbol{\beta}_1 - t\boldsymbol{v}_1\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} \\ &+ \int_{S_2} \|\boldsymbol{x} - \boldsymbol{\beta}_2 - t\boldsymbol{v}_2\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} + \text{constant} \\ &= \int_{S_1} \|\boldsymbol{x} - \boldsymbol{s}_1\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} + \int_{S_1} \|\boldsymbol{s}_1 - \boldsymbol{\beta}_1 - t\boldsymbol{v}_1\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} \\ &+ \int_{S_2} \|\boldsymbol{x} - \boldsymbol{s}_2\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} + \int_{S_2} \|\boldsymbol{s}_2 - \boldsymbol{\beta}_2 - t\boldsymbol{v}_2\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} \\ &+ \text{constant.} \end{split}$$

(In the calculation above we have avoided double counting the contribution from $\mathcal{V}_1(\beta) = \mathcal{V}_2(\beta)$, and the constant part corresponds to the contribution from $\mathcal{V}_i(\beta)$ for $i \geq 3$) With the above choices of v_1 and v_2 , we see that $h^{v}(0) > h^{v}(t)$ for all $t \in (0, 1)$ and hence t = 0 is not a local minimum of h^{v} , which is a contradiction. Therefore, we must have $\beta_1 \neq \beta_2$ whenever $\mathcal{V}_i(\beta) \cup \mathcal{V}_j(\beta)$ has a positive measure. The more general statement in Lemma 2 can be established in a similar manner.

Now suppose that β_i has a Voronoi set $\mathcal{V}_i(\beta)$ with a positive measure. In this case the center of mass $c_i := \frac{\int_{\mathcal{V}_i(\beta)} xf(x)dx}{\int_{\mathcal{V}_i(\beta)} f(x)dx}$ is well-defined. Choose the direction $v = (0, \dots, 0, c_i - \beta_i, 0, \dots, 0)$. Since t = 0 is a local minimum of H^v , its derivative must vanish at t = 0. Using the derivative expression from Lemma $[1]^4$ we obtain that

$$egin{aligned} 0 &= rac{\mathrm{d}}{\mathrm{d}t} H^{m{v}}(0) = -\int_{\mathcal{V}_i(m{eta})} 2 \langle m{v}_i,m{x} - m{eta}_i
angle f(m{x}) \mathrm{d}m{x} \ &= -2 \langle m{v}_i,m{c}_i - m{eta}_i
angle \int_{\mathcal{V}_i(m{eta})} f(m{x}) \mathrm{d}m{x} \ &= -2 \|m{c}_i - m{eta}_i\|^2 \int_{\mathcal{V}_i(m{eta})} f(m{x}) \mathrm{d}m{x}, \end{aligned}$$

where the last step follows from our choice of v. Since $\int_{\mathcal{V}_i(\beta)} f(x) dx$ is the measure of $\mathcal{V}_i(\beta)$ and positive, we must have $\beta_i = c_i$ as claimed.

We state and prove several technical lemmas that are used in Section $\boxed{\nabla I}$

Proof of Lemma 4

Recall that $m_{i,s}$ and $c_{i,s}$ denote the mass and the center of mass of the set \mathcal{V}_i with respect to the density f_s . We similarly define

$$\begin{split} \widetilde{m_{i,s}} &= \sum_{s' \in T_i: s' \neq s} m_{i,s'} \quad \text{and} \\ \widetilde{c_{i,s}} &= \frac{\sum_{s' \in T_i: s' \neq s} m_{i,s'} c_{i,s'}}{\sum_{s' \in T_i: s' \neq s} m_{i,s'}}, \end{split}$$

which are the mass and the center of mass of the set \mathcal{V}_i with respect to the density $\sum_{s'\neq s} f_{s'}$. With this notation, the local minimum β must satisfy the necessary condition

$$\boldsymbol{\beta}_{i} = \frac{m_{i,s} \boldsymbol{c}_{i,s} + \widetilde{m_{i,s}} \boldsymbol{c}_{i,s}}{m_{i,s} + \widetilde{m_{i,s}}}, \tag{36}$$

which follows from Lemma 2 and the text thereafter. Rearranging the expression (36) gives

$$\widetilde{\boldsymbol{c}_{i,s}} = \frac{(\widetilde{m_{i,s}} + m_{i,s})\boldsymbol{\beta}_i - m_{i,s}\boldsymbol{c}_{i,s}}{\widetilde{m_{i,s}}},$$

It then follows from the triangle inequality that

$$\begin{aligned} \|\widetilde{\boldsymbol{c}_{i,s}} - \boldsymbol{\beta}_{i}\| &= \frac{m_{i,s}}{\widetilde{m_{i,s}}} \|\boldsymbol{c}_{i,s} - \boldsymbol{\beta}_{i}\| \\ &\leq \frac{m_{i,s}}{\widetilde{m_{i,s}}} (\|\boldsymbol{c}_{i,s} - \boldsymbol{\beta}_{s}^{*}\| + \|\boldsymbol{\beta}_{s}^{*} - \boldsymbol{\beta}_{i}\|). \end{aligned} (37)$$

⁴Lemma \mathbf{I} is applicable for the following reason: we can ignore those $\mathcal{V}_i(\beta)$'s with zero measure in the integrals defining G and H^v , in which case we have just established that β must have pairwise distinct components and thus satisfy the premise of Lemma \mathbf{I} .

We are now ready to prove Lemma 4, whose assumption states that $\rho_s(\partial_{j,\ell}) > \lambda = \frac{c}{\sqrt{r\Delta_{\max}}}$ for some $(s, j, \ell) \in T_i \times [k] \times [k]$. Observation 1] ensures that such an s is unique, hence for all other $s' \in T_i \setminus \{s\}$, we must have $\rho_{s'}(\partial_{j,\ell}) \leq \lambda, \forall (j, \ell)$. Observation 2] ensures that for all these s', if $\beta_{s'}^* \in \mathcal{V}_i$ then $s' \in A_i$. In view of these properties and equation (36), we can see that $\widetilde{c_{i,s}}$ is similar to β_i except that the density of the s-th true cluster is ignored. Therefore, we can follow the same arguments for proving Part 2(b) of Theorem 4 as well as the simplification step in (12) to obtain that

$$\|\widetilde{\boldsymbol{c}_{i,s}} - \boldsymbol{b}_i^-\| \le \Delta_{\max} \frac{(4k^* + 2ck^2)(1 + \frac{2k^*}{k})}{\sqrt{\eta_{\max}}} \qquad (38)$$

On the other hand, under the assumption of the lemma, Part 1 of Theorem 4 ensures that

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{s}^{*}\| \leq \frac{k}{\lambda} + 3r \leq \Delta_{\max} \frac{4k^{*}}{c\sqrt{\eta_{\max}}}.$$
 (39)

We proceed by bounding $\|\boldsymbol{b}_i^- - \boldsymbol{\beta}_s^*\|$ as follows:

$$\begin{aligned} \|\boldsymbol{b}_{i}^{-} - \boldsymbol{\beta}_{s}^{*}\| \leq \|\boldsymbol{b}_{i}^{-} - \widetilde{\boldsymbol{c}_{i,s}}\| + \|\widetilde{\boldsymbol{c}_{i,s}} - \boldsymbol{\beta}_{i}\| + \|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{s}^{*}\| \\ \leq \|\boldsymbol{b}_{i}^{-} - \widetilde{\boldsymbol{c}_{i,s}}\| + \frac{m_{i,s}}{\widetilde{m}_{i,s}}(\|\boldsymbol{c}_{i,s} - \boldsymbol{\beta}_{s}^{*}\| + \|\boldsymbol{\beta}_{s}^{*} - \boldsymbol{\beta}_{i}\|) \\ + \|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}^{*}\| \end{aligned}$$

$$(40)$$

$$\leq \Delta_{\max} \frac{(4k^* + 2ck^2)(1 + \frac{2k^*}{k})}{\sqrt{\eta_{\max}}} + \frac{m_{i,s}}{\widetilde{m_{i,s}}} \cdot \frac{\Delta_{\max}}{\eta_{\max}} + \left(\frac{m_{i,s}}{\widetilde{m_{i,s}}} + 1\right) \Delta_{\max} \frac{4k^*}{c\sqrt{\eta_{\max}}},$$
(41)

where step (40) follows from the bound (37), and step (41) follows from the bounds (38) and (39) as well as the fact that $c_{i,s} \in \mathbb{B}_s$ so $||c_{i,s} - \beta_s^*|| \leq r = \frac{\Delta_{\max}}{\eta_{\max}}$. Now, note that since $|A_i \setminus \{s\}| \geq 2$ by assumption, there exists some $s' \in A_i \subseteq T_i$ such that $s' \neq s$. We have established above that this s' must satisfy $\rho_{s'}(\partial_{j,\ell}) \leq \lambda, \forall (j,\ell)$, hence applying Part 2 of Theorem 4 we obtain that $m_{i,s'} = \mathbb{P}_{s'}(\mathcal{V}_i) \geq 0.5$, which further implies $\frac{m_{i,s}}{m_{i,s}} \leq 2$. Continuing from the above display equation, we obtain

$$\begin{aligned} \|\boldsymbol{b}_{i}^{-} - \boldsymbol{\beta}_{s}^{*}\| \leq & \Delta_{\max} \frac{(4k^{*} + 2ck^{2})(1 + \frac{2k^{*}}{k})}{\sqrt{\eta_{\max}}} \\ &+ 3\Delta_{\max} \frac{4k^{*}}{c\sqrt{\eta_{\max}}} + \Delta_{\max} \frac{2}{\eta_{\max}} \\ \leq & \Delta_{\max} \frac{2ck^{2}(1 + \frac{2k^{*}}{k}) + 4k^{*}(2.5 + \frac{2k^{*}}{k})}{\sqrt{\eta_{\max}}}, \end{aligned}$$
(42)

where the last step follows from the assumption that c > 3and $\eta_{\text{max}} \ge 4ck^2$. Combining the inequalities (39) and (42), we obtain

$$\begin{aligned} \|\boldsymbol{b}_{i}^{-} - \boldsymbol{\beta}_{i}\| \leq \|\boldsymbol{b}_{i}^{-} - \boldsymbol{\beta}_{s}^{*}\| + \|\boldsymbol{\beta}_{s}^{*} - \boldsymbol{\beta}_{i}\| \\ \leq \Delta_{\max} \frac{2ck^{2}(1 + \frac{2k^{*}}{k}) + 4k^{*}(3 + \frac{2k^{*}}{k})}{\sqrt{\eta_{\max}}}, \quad (43) \end{aligned}$$

thereby proving the first bound in Lemma 4.

To prove the second bound in Lemma 4, we observe that by definition of b_i^+ and b_i^- , there holds

$$\begin{split} \boldsymbol{b}_{i}^{+} &:= \frac{1}{1 + |A_{i} \setminus \{s\}|} \left(\sum_{s' \in A_{i} \setminus \{s\}} \beta_{s'}^{*} + \beta_{s}^{*} \right) \\ &= \frac{|A_{i} \setminus \{s\}|}{1 + |A_{i} \setminus \{s\}|} \boldsymbol{b}_{i}^{-} + \frac{1}{1 + |A_{i} \setminus \{s\}|} \beta_{s}^{*}. \end{split}$$

whence $\|\boldsymbol{b}_i^+ - \boldsymbol{\beta}_s^*\| = \frac{|A_i \setminus \{s\}|}{1 + |A_i \setminus \{s\}|} \|\boldsymbol{b}_i^- - \boldsymbol{\beta}_s^*\|$. It follows that

$$\begin{split} \|\boldsymbol{b}_{i}^{+} - \boldsymbol{\beta}_{i}\| \leq & \|\boldsymbol{b}_{i}^{+} - \boldsymbol{\beta}_{s}^{*}\| + \|\boldsymbol{\beta}_{s}^{*} - \boldsymbol{\beta}_{i}\| \\ \leq & \|\boldsymbol{b}_{i}^{-} - \boldsymbol{\beta}_{s}^{*}\| + \|\boldsymbol{\beta}_{s}^{*} - \boldsymbol{\beta}_{i}\| \\ \leq & \Delta_{\max} \frac{2ck^{2}(1 + \frac{2k^{*}}{k}) + 4k^{*}(3 + \frac{2k^{*}}{k})}{\sqrt{\eta_{\max}}} \end{split}$$

where the last step follows from equation (43).

It remains to show that $|A_i \setminus \{s\}| \ge 2$. Note that $A_i \ne \emptyset$ under the assumption $|W_i \setminus \{s\}| \ge 1$ of the lemma. For the sake of deriving a contradiction, assume that $W_i \setminus \{s\} = \{s'\}$, in which case $\boldsymbol{b}_i^- = \boldsymbol{\beta}_{s'}^*$. It then follows from inequality (42) that $\|\boldsymbol{\beta}_{s'}^* - \boldsymbol{\beta}_s^*\| \le \Delta_{\max} \frac{2ck^2(1+\frac{2k^*}{k})+4k^*(3+\frac{2k^*}{k})}{\sqrt{\eta_{\max}}}$, contradicting the separation assumption on η_{\min} in Theorem 3. This completes the proof of Lemma 4.

Controlling the Volume

In this section, we show that the intersection of a Voronoi set and a ground truth cluster must be small if (i) the true center is not in the Voronoi set and (ii) the intersection of the true cluster and the boundary of the Voronoi set is small. This is formalized in the following lemma.

Lemma 8 (Controlling the volume of intersection). Let μ be the uniform distribution on $\mathbb{B}_{\mathbf{0}}(r)$. Let P be a closed polyhedron with at most k facets satisfying $\mathbf{0} \notin \operatorname{int}(P)$. If each facet F of P satisfies $\frac{1}{r^d V_d} \operatorname{ReVol}(F \cap \mathbb{B}_{\mathbf{0}}(r)) \leq \lambda$, then we have $\mu(P) \leq k \lambda r$.

Proof. Introduce the shorthand $\mathbb{B} := \mathbb{B}_{\mathbf{0}}(r)$. We may assume that $P \cap \mathbb{B} \neq \emptyset$, because otherwise the lemma is trivially true. We claim that one may shift the polyhedron P by a distance r so that its intersection with the ball \mathbb{B} has zero measure. That is, there exists a unit vector v such that $(P+(r+\epsilon)v) \cap \mathbb{B} = \emptyset$ for all $\epsilon > 0$. We further claim that v can be chosen in such a way that the intersection $P \cap \mathbb{B}$ is enclosed by the original boundary and the shifted boundary; that is, $P \cap \mathbb{B} \subseteq (\partial P \cap \mathbb{B}) + L_r$, where we $L_r := \{tv : t \in [0, r]\}$ is a line segment. Figure [0, r] provides an illustration of these two claims, whose proof is deferred to the end of this section. Therefore, we have the bound

$$\begin{split} \mu(P) = & \frac{\operatorname{Vol}(P \cap \mathbb{B})}{\operatorname{Vol}(\mathbb{B})} \leq \frac{\operatorname{Vol}((\partial P \cap \mathbb{B}) + L_r)}{\operatorname{Vol}(\mathbb{B})} \\ \leq & \frac{r \sum_{F \in \mathcal{F}} \operatorname{Re} \operatorname{Vol}(F \cap \mathbb{B})}{r^d V_d} \leq rk\lambda, \end{split}$$

where $\mathcal{F} := \{F : F \text{ is a facet of } P\}$ satisfies $|\mathcal{F}| \leq k$ by assumption. This completes the proof of the lemma.

Let us prove the two claims above. Because P is convex and $\mathbf{0} \notin \operatorname{int}(P)$, the separating hyperplane theorem ensures that



Figure 6. Shifting the boundary of the polyhedron P to bound the volume of $P \cap \mathbb{B}_{\mathbf{0}}(r)$.

there exists some unit vector v such that $\langle x, v \rangle \ge 0, \forall x \in P$. Therefore, for all $x \in P$ and $\epsilon > 0$, we have

$$\begin{aligned} \|\boldsymbol{x} + (r+\epsilon)\boldsymbol{v}\|^2 &= \|\boldsymbol{x}\|^2 + 2(r+\epsilon)\langle \boldsymbol{x}, \boldsymbol{v} \rangle + (r+\epsilon)^2 \|\boldsymbol{v}\|^2 \\ &\geq 0 + 0 + (r+\epsilon)^2 > r^2, \end{aligned}$$

whence $\mathbf{x} + (r + \epsilon)\mathbf{v} \notin \mathbb{B}$, proving the first claim. To prove the second claim, fix an arbitrary $\mathbf{x} \in P \cap \mathbb{B}$ and consider the half line $\ell := \{\mathbf{x} - t\mathbf{v} : t \ge 0\}$. Note that ℓ must intersect the boundary ∂P ; otherwise we would have $\ell \subseteq P$ and hence the separating hyperplane property implies that $\langle \mathbf{x} - t\mathbf{v}, \mathbf{v} \rangle \ge 0$ for all $t \ge 0$, which cannot hold as \mathbf{v} has unit norm. Since P is convex, ℓ intersects ∂P at a unique point, say $\mathbf{x}_0 =$ $\mathbf{x} - t_0 \mathbf{v}$. We must have $t_0 \le r$; otherwise we would have $\mathbf{x} = \mathbf{x}_0 + t_0 \mathbf{v} \in P + (r + \epsilon)\mathbf{v}$ for some $\epsilon > 0$ and hence $\mathbf{x} \notin \mathbb{B}$ by the first claim, which is a contradiction. Using the separating hyperplane property $0 \le \langle \mathbf{x} - t_0 \mathbf{v}, \mathbf{v} \rangle \le \langle \mathbf{x}, \mathbf{v} \rangle$ again, we have

$$\begin{aligned} \|\boldsymbol{x}_0\|^2 &= \|\boldsymbol{x} - t_0 \boldsymbol{v}\|^2 \\ &= \|\boldsymbol{x}\|^2 + \langle -t_0 \boldsymbol{v}, \boldsymbol{x} \rangle + \langle \boldsymbol{x} - t_0 \boldsymbol{v}, -t_0 \boldsymbol{v} \rangle \le r^2 + 0 + 0 \end{aligned}$$

and thus $\boldsymbol{x}_0 \in \mathbb{B} \cap \partial P$. Combining pieces, we conclude that $\boldsymbol{x} = \boldsymbol{x}_0 + t_0 \boldsymbol{v} \in (\partial P \cap \mathbb{B}) + L_r$. As $\boldsymbol{x} \in P \cap \mathbb{B}$ is arbitrary, we have $P \cap \mathbb{B} \subseteq (\partial P \cap \mathbb{B}) + L_r$ as claimed. \Box

Below we state a analogous version of Lemma 8 under the Gaussian density. This result is used in the proof of our main Theorem 5 for the Gaussian Mixture model.

Lemma 9 (Controlling the volume of intersection, Gaussian). Let μ be the probability measure with respect to a spherical Gaussian distribution in \mathbb{R}^d with mean **0** and variance σ^2 . Let P be a closed polyhedron with at most k facets satisfying **0** \notin int(P). If each facet F of P satisfies $\rho_1(F) \leq \lambda$, where $\rho_1(F)$ is defined as in equation (48) with $\beta_1^* = \mathbf{0}$, then we have $\mu(P \cap \mathbb{B}_0(r)) \leq k\lambda r$ for any r > 0.

Proof. We fix r and introduce the shorthand $\mathbb{B} \equiv \mathbb{B}_{0}(r)$. From the proof of Lemma 8, there exists a unit vector $\boldsymbol{v} = (v_1, \ldots, v_d)$ satisfying (i) $\langle \boldsymbol{v}, \boldsymbol{x} \rangle \geq 0$ for all $\boldsymbol{x} \in P$; (ii) $P \cap \mathbb{B} \subseteq \partial P \cap \mathbb{B} + L_r$, where $L_r := \{t\boldsymbol{v} : t \in [0, r]\}$. It follows that

$$\mu(P \cap \mathbb{B}) \le \mu(\partial P \cap \mathbb{B} + L_r) \le \sum_{F \in \mathcal{F}} \mu(F \cap \mathbb{B} + L_r),$$

where $\mathcal{F} := \{F : F \text{ is a facet of } P\}$. Since by assumption $|\mathcal{F}| \leq k$ and each facet $F \in \mathcal{F}$ satisfies $\rho_1(F) \leq \lambda$, the lemma follows if we can show that $\mu(F \cap \mathbb{B} + L_r) \leq \rho_1(F) \cdot r$.

Without loss of generality, we may use an orthonormal basis $\{w_1, \ldots, w_d\}$ such that w_1 is the normal vector of F and $\langle v, w_1 \rangle \geq 0$. Under this basis, we have $x_1 = \langle x, w_1 \rangle = c$ for all $x \in F$, where c is a fixed number. We define a vertical slice at $x_1 = c$ as follows:

$$S_F := \{(x_2, \ldots, x_d) : (c, x_2, \ldots, x_d) \in F \cap \mathbb{B}\}.$$

With this notation, the definition of $\rho_1(F)$ is equivalent to

 $\rho_1(F)$

$$= \int_{(x_2,...,x_d)\in S_F} \frac{1}{(\sqrt{2\pi\sigma^2})^d} \exp\left(-\frac{c^2 + \sum_{j=2}^d x_j^2}{2\sigma^2}\right) \mathrm{d}x_2 \dots \mathrm{d}x_d,$$

and the set $F \cap \mathbb{B} + L_r$ admits the explicit expression

$$F \cap \mathbb{B} + L_r = \{ \boldsymbol{z} : \boldsymbol{z} = (c, x_2, \dots, x_d) + t\boldsymbol{v}, \\ (x_2, \dots, x_d) \in S_F, t \le r \}.$$
(44)

We are now ready to bound the quantity of interest

$$\mu(F \cap \mathbb{B} + L_r) = \int_{\boldsymbol{z} \in F \cap \mathbb{B} + L_r} \frac{1}{(\sqrt{2\pi\sigma^2})^d} \exp\left(-\frac{\sum_{j=1}^d z_j^2}{2\sigma^2}\right) \mathrm{d}z_1 \dots \mathrm{d}z_d$$

Using the expression (44), we may rewrite the above integral by a change of variable from (z_1, z_2, \ldots, z_d) to (t, x_2, \ldots, x_d) . To this end, note that the corresponding Jacobian matrix is

$$J = \begin{bmatrix} v_1 & 0 & 0 & \dots & 0 \\ v_2 & 1 & 0 & \dots & 0 \\ \vdots & & & & \\ v_d & 0 & 0 & \dots & 1 \end{bmatrix},$$

which is a lower triangular matrix satisfying $|\det(J)| = |v_1|$. Consequently, we have

$$\mu(F \cap \mathbb{B} + L_r)$$

$$= \int_0^r \left[\int_{(x_2, \dots, x_d) \in S_F} \frac{1}{(\sqrt{2\pi\sigma^2})^d} \cdot \exp\left(-\frac{(c+tv_1)^2 + \sum_{j=2}^d (x_j + tv_j)^2}{2\sigma^2}\right) \mathrm{d}x_2 \dots \mathrm{d}x_d \right] |v_1| \mathrm{d}t.$$

To upper bound the Gaussian density above, we note that for any $t \ge 0$ and $(c, x_2, \ldots, x_d) \in F \subseteq P$, the property of the separating hyperplane ensures that $t(v_1c + \sum_{j=2}^d v_j x_j) \ge 0$, whence

$$\exp\left(-\frac{(c+tv_1)^2 + \sum_{j=2}^d (x_j+tv_j)^2}{2\sigma^2}\right)$$
$$= \exp\left(-\frac{c^2 + \sum_{j=2}^d (x_j)^2}{2\sigma^2}\right) \exp\left(-\frac{t(v_1c + \sum_{j=2}^d v_j x_j)}{\sigma^2}\right) \cdot \exp\left(-\frac{t^2}{2\sigma^2}\right)$$
$$\leq \exp\left(-\frac{c^2 + \sum_{j=2}^d x_j^2}{2\sigma^2}\right).$$

It follows that

$$\mu(F \cap \mathbb{B} + L_r) \leq \int_0^r \left[\int_{(x_2, \dots, x_d) \in S_F} \frac{1}{(\sqrt{2\pi\sigma^2})^d} \right]$$

$$\exp\left(-\frac{c^2 + \sum_{j=2}^d (x_j)^2}{2\sigma^2}\right) \mathrm{d}x_2 \dots \mathrm{d}x_d \Big] |v_1| \mathrm{d}t$$
$$= \int_0^r \rho_1(F) \cdot |v_1| \mathrm{d}t \le \rho_1(F) \cdot r,$$

where the last step holds since \boldsymbol{v} is a unit vector with $|v_1| \leq 1$. This completes the proof of Lemma 9.

Controlling the distance to the center

In this section, we prove the following result:

Lemma 10 (Bound on the center of mass, ball). Let μ be the uniform distribution over the ball $\mathbb{B}_0(r) \subset \mathbb{R}^d$. Suppose that a subset $S \subset \mathbb{R}^d$ has probability measure $\mu(S) > 0$. Let c_S be the center of mass of the set S with respect to μ . We have the bound

$$\|\boldsymbol{c}_S\| \le \frac{r \cdot \mu(\mathbb{R}^d \setminus S)}{\mu(S)}$$

Proof. Recall the expression of the center of mass $c_S = \frac{\int x \mathbf{1}_S(x) d\mu}{\mu(S)}$. We have

$$\begin{split} \mu(S) \cdot \|\boldsymbol{c}_S\| &= \left\| \int \boldsymbol{x} \mathbb{1}_S(\boldsymbol{x}) \mathrm{d} \mu \right\| \\ &\stackrel{(i)}{=} \left\| - \int \boldsymbol{x} \mathbb{1}_{\mathbb{R}^d \setminus S}(\boldsymbol{x}) \mathrm{d} \mu \right\| \\ &\stackrel{(ii)}{\leq} \int \|\boldsymbol{x}\| \mathbb{1}_{\mathbb{R}^d \setminus S}(\boldsymbol{x}) \mathrm{d} \mu \\ &\stackrel{(iii)}{\leq} r \mu(\mathbb{R}^d \setminus S), \end{split}$$

where step (i) holds since μ has mean **0**, and step (ii) holds by the Jensen's inequality, and step (iii) holds because $||\mathbf{x}|| \leq r$ for all $\mathbf{x} \in \text{support}(\mu) = \mathbb{B}_{\mathbf{0}}(r)$. Rearranging the inequality proves the desired bound.

Lemma 11 (Bound on the center of mass, Gaussian). Let μ be the Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$. Suppose that a subset $S \subset \mathbb{R}^d$ has probability measure $\mu(S) > 0$. Let \mathbf{c}_S be the center of mass of the set S with respect to μ . We have the bound

$$\|\boldsymbol{c}_S\| \le \frac{2\sigma \cdot \sqrt{\mu(\mathbb{R}^d \setminus S)}}{\mu(S)}$$

Proof. Recall the variational characterization of the center of mass c_S :

Since **0** is the mean of μ , we have

$$0 \leq \|\boldsymbol{c}_{S}\|^{2}$$

$$= \int \|\boldsymbol{x} - \boldsymbol{c}_{S}\|^{2} d\mu - \int \|\boldsymbol{x}\|^{2} d\mu$$

$$= \int \|\boldsymbol{x} - \boldsymbol{c}_{S}\|^{2} \mathbb{1}_{S}(\boldsymbol{x}) d\mu - \int \|\boldsymbol{x}\|^{2} \mathbb{1}_{S}(\boldsymbol{x}) d\mu$$

$$+ \int \|\boldsymbol{x} - \boldsymbol{c}_{S}\|^{2} \mathbb{1}_{\mathbb{R}^{d} \setminus S}(\boldsymbol{x}) d\mu - \int \|\boldsymbol{x}\|^{2} \mathbb{1}_{\mathbb{R}^{d} \setminus S}(\boldsymbol{x}) d\mu$$

$$\stackrel{(i)}{\leq} \int \|\boldsymbol{x} - \boldsymbol{c}_{S}\|^{2} \mathbb{1}_{\mathbb{R}^{d} \setminus S}(\boldsymbol{x}) d\mu - \int \|\boldsymbol{x}\|^{2} \mathbb{1}_{\mathbb{R}^{d} \setminus S}(\boldsymbol{x}) d\mu$$

$$= \mu(\mathbb{R}^{d} \setminus S) \|\boldsymbol{c}_{S}\|^{2} - 2 \int \langle \boldsymbol{x}, \boldsymbol{c}_{S} \rangle \mathbb{1}_{\mathbb{R}^{d} \setminus S}(\boldsymbol{x}) d\mu. \quad (45)$$

where step (i) follows from the variational characterization of the center of mass

$$oldsymbol{c}_S = \operatorname{argmin}_{oldsymbol{z} \in \mathbb{R}^d} \int \|oldsymbol{x} - oldsymbol{z}\|^2 \mathbbm{1}_S(oldsymbol{x}) \mathrm{d} \mu.$$

Rearranging equation (45) gives

$$\begin{split} \mu(S) \|\boldsymbol{c}_S\|^2 &\leq -2 \int \langle \boldsymbol{x}, \boldsymbol{c}_S \rangle \mathbb{1}_{\mathbb{R}^d \setminus S}(\boldsymbol{x}) \mathrm{d}\mu \\ &\stackrel{\text{(ii)}}{\leq} 2 \|\boldsymbol{c}_S\| \sqrt{\int \left\langle \boldsymbol{x}, \frac{\boldsymbol{c}_S}{\|\boldsymbol{c}_S\|} \right\rangle^2} \mathrm{d}\mu \cdot \sqrt{\int \mathbb{1}_{\mathbb{R}^d \setminus S}(\boldsymbol{x}) \mathrm{d}\mu} \\ &\stackrel{\text{(iii)}}{=} 2 \|\boldsymbol{c}_S\| \sigma \cdot \sqrt{\mu(\mathbb{R}^d \setminus S)}, \end{split}$$

where step (ii) follows from Cauchy-Schwarz and step (iii) follows from the fact that any one-dimensional margin of $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$ is the univariate Gaussian distribution $\mathcal{N}(0, \sigma^2)$. Rearranging the above equation proves the desired bound. \Box

PROOF OF THEOREM 2

We prove a more general and quantitative version Theorem 2 that holds for any k^* and k.

Theorem 5 (General version of Theorem 2). Let t be any number satisfying $\varphi(t) := 2 \exp(-t^2 \min(d, k + k^*)/8) < \frac{1}{4}$. Under the Gaussian Mixture Model, assume that $\eta_{\max} \ge 16tc^2k^4$ and $\eta_{\min} \ge \left[\frac{k^*}{k\sqrt{t}}\left(2 + \frac{c\sqrt{t}}{\sqrt{\eta_{\max}}}\right) + 3ck(k^* + k)\sqrt{t}\right]\sqrt{\eta_{\max}} + 2(2k^* + 1)\varphi(t)\eta_{\max}$ for some constant $c \ge 3$. If $\beta = (\beta_1, \ldots, \beta_k) \in \mathbb{R}^{d \times k}$ is a local minimum of G, then the ground truth centers and fitted centers can be partitioned as $[k] = \bigcup_{a=1}^m S_a^*$ and $[k] = \bigcup_{a=0}^m S_a$, respectively, such that for each $a \in [m]$, exactly one of the following holds:

• (many/one-fit-one association) $|S_a| \ge 1$ and $S_a^* = \{s\}$ for some $s \in [k]$; moreover,

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{\beta}_{s}^{*}\| \leq \Delta_{\max} \left\{ \frac{\frac{2k^{*}}{ck\sqrt{t}} + 2ck(k^{*} + k)\sqrt{t}}{\sqrt{\eta_{\max}}} + 7k^{*}\varphi(t) \right\}$$

$$\forall i \in S_{a}.$$
 (46)

• (one-fit-many association) $S_a = \{i\}$ for some $i \in [k]$ and $|S_a^*| \ge 2$; moreover,

$$\left\|\boldsymbol{\beta}_{i} - \frac{1}{|S_{a}^{*}|} \sum_{s \in S_{a}^{*}} \boldsymbol{\beta}_{s}^{*}\right\|$$

$$\leq \Delta_{\max} \left\{ \frac{\frac{2k^{*}}{ck\sqrt{t}} + 3ck(k^{*} + k)\sqrt{t}}{\sqrt{\eta_{\max}}} + 7k^{*}\varphi(t) \right\}$$
(47)

In addition, for each $i \in S_0$, we have $\mathbb{P}(\mathcal{V}_i(\beta)) \leq \frac{ck\sqrt{t}}{\sqrt{\eta_{\max}}} + \varphi(t)$ (almost-empty association).

Given Theorem 5. Theorem 2 follows immediately. The rest of this section is devoted to proving Theorem 5.

As before, we use \mathbb{P} to denote the probability measure with respect to f and \mathbb{P}_s to denote the probability measure with respect to f_s , where f is the density of the Gaussian mixture and f_s is the density of the *s*-th Gaussian component. Recall the population *k*-mean objective function:

$$G(\boldsymbol{\beta}) = \int_{\boldsymbol{x}} \min_{j \in [k]} \|\boldsymbol{x} - \boldsymbol{\beta}_j\|^2 f(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}$$

$$= \frac{1}{k^*} \sum_{s=1}^{k^*} \int_{\boldsymbol{x}} \min_{j \in [k]} \|\boldsymbol{x} - \boldsymbol{\beta}_j\|^2 f_s(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}.$$

Reduction to lower dimensions

We first argue that it suffices to prove the theorem in dimension $d' \leq k + k^*$. Once this is established, then the theorem for $d > k^* + k$ dimensions can be deduced as follows. Suppose that $\beta^* \in \mathbb{R}^{d \times k^*}$ is the ground truth solution and $\beta \in \mathbb{R}^{d \times k}$ is a candidate solution. We may choose a coordinate system such that the first $d' = k^* + k$ dimensions correspond to $\operatorname{span} \{\beta_1^*, \ldots, \beta_{k^*}^*, \beta_1, \ldots, \beta_k\}$. In this case, for each $s \in [k^*]$ we have $\beta_s^* = (\beta_s'', \mathbf{0})$ for some $\beta_s' \in \mathbb{R}^{d'}$ and for each $i \in [k]$, we have and $\beta_i = (\beta_i', \mathbf{0})$ for some $\beta_i' \in \mathbb{R}^{d'}$. Moreover, thanks to Gaussian's rotational invariance, the *d*-dimensional Gaussian mixture is a product distribution with respect to the first *d'* dimensions and the last d - d' dimensions, where the first *d'*-dimensional margin is itself a Gaussian mixture. Indeed, for any $\mathbf{x} = (\mathbf{x}', \mathbf{z}) \in \mathbb{R}^d$ with $\mathbf{x}' \in \mathbb{R}^{d'}$ and $\mathbf{z} \in \mathbb{R}^{d-d'}$, the density of the Gaussian mixture factorizes:

$$\begin{split} f(\boldsymbol{x}) &\propto \frac{1}{k^*} \sum_{s=1}^{k^*} \exp\left(\frac{\|\boldsymbol{x} - \boldsymbol{\beta}_s^*\|^2}{2\sigma^2}\right) \\ &= \frac{1}{k^*} \sum_{s=1}^{k^*} \exp\left(\frac{\|(\boldsymbol{x}', \boldsymbol{z}) - (\boldsymbol{\beta}_s^{*\prime}, \boldsymbol{0})\|^2}{2\sigma^2}\right) \\ &= \left[\frac{1}{k^*} \sum_{s=1}^{k^*} \exp\left(\frac{\|\boldsymbol{x}' - \boldsymbol{\beta}_s^{\prime\prime}\|^2}{2\sigma^2}\right)\right] \cdot \exp\left(\frac{\|\boldsymbol{z}\|^2}{2\sigma^2}\right). \end{split}$$

Now, if β is a local minimum of G, then β' is also a local minimum of G restricted to the first $d' = k^* + k$ dimensions. Applying the theorem with dimension d', we obtain bounds on the quantities $\|\beta'_i - \beta^{*'}_s\|$, $\|\beta'_i - \sum_{s \in S} \beta^{*'}_s\|$ and $\mathbb{P}(\mathcal{V}_i(\beta'))$. We claim that these three quantities are equal to $\|\beta_i - \beta^*_s\|$, $\|\beta_i - \sum_{s \in S} \beta^*_s\|$ and $\mathbb{P}(\mathcal{V}_i(\beta))$, respectively. Indeed, the first two equalities are immediate under our coordinate system; the last equality holds because the Gaussian mixture factorzes (shown above) and so do the Voronoi sets: $\mathcal{V}_i(\beta) = \mathcal{V}_i(\beta') \times \mathbb{R}^{d-d'}$. We conclude that the same collection of bounds hold in dimension d as well. In the rest of the proof, we can safely assume that $d \leq k^* + k$, in which case $\min\{k^* + k, d\} = d$.

As in the proof for the Stochastic Ball Model, we first establish a general result, an analogue of Theorem 4 that provides a family of bounds parametrized by two positive numbers λ and t. To state this result, we introduce some additional notation. Set $r = t\sigma\sqrt{d}$ and let

$$\mathbb{B}_s(r) := \{ \boldsymbol{x} \in \mathbb{R}^d : \| \boldsymbol{x} - \boldsymbol{\beta}_s^* \| \le r \}.$$

denote the ball centered at β_s^* with radius r. Recall that for each $i, j \in [k]$ and $s \in [k^*]$, the Voronoi boundary $\partial_{i,j}$ lies in a (d-1)-dimensional affine subspace \mathcal{L} . Since the distribution f_s of the s-th component is rotationally invariant, we may assume WLOG that $\mathcal{L} = \{x \in \mathbb{R}^d : x_1 = z\}$ for some number z. Accordingly, we define the quantity

$$\rho_s(\partial_{i,j}) := \int \mathbb{1}\left\{ (z, x_2, \dots, x_d) \in \partial_{i,j} \cap \mathbb{B}_s(r) \right\}.$$

$$f_s(z, x_2, \dots, x_d) \,\mathrm{d}x_2 \dots \mathrm{d}x_d,\tag{48}$$

which is a measure of the relative probability mass of the Voronoi boundary $\partial_{i,j}$ when restricted to the ball $\mathbb{B}_s(r)$. Also recall the function $\varphi(\cdot)$ defined in the statement of Theorem which satisfies $\varphi(t) = 2 \exp(-t^2 d/8)$ when $d \leq k^* + k$.

We now state the family of bound in the Gaussian case:

Theorem 6 (Family of bounds for Gaussian). Under the Gaussian mixture model, let $\beta = (\beta_1, ..., \beta_k)$ be a local minima for the k-means objective function G defined in (3). Let $\lambda > 0$ and t > 0 be two arbitrary fixed numbers and set $r := t\sigma\sqrt{d}$. For each $i \in [k]$, define the sets:

$$T_i := \{ s \in [k^*] : \mathcal{V}_i \cap \mathbb{B}_s(r) \neq \emptyset \} \text{ and } A_i := \{ s \in [k^*] : \beta_s^* \in \operatorname{int}(\mathcal{V}_i) \} \subseteq T_i.$$

Then the following is true for each $i \in [k]$:

1) If $\rho_s(\partial_{j,\ell}) > \lambda$ for some pair (j,ℓ) and $s \in T_i$, then

$$\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_s^*\| \le \frac{k^*}{\lambda} + 3r.$$

2) For each $s \in T_i$, if $\rho_s(\partial_{j,\ell}) \leq \lambda$ for all pairs (j, ℓ) , then the following bounds hold

$$\mathbb{P}_{s}(\mathcal{V}_{i}) \geq 1 - k^{2}\lambda r - \varphi(t), \quad \forall s \in A_{i}, \\ \mathbb{P}_{s}(\mathcal{V}_{i}) \leq k\lambda + \varphi(t), \qquad \forall s \in T_{i} \setminus A_{i}.$$

Furthermore, if $\rho_s(\partial_{j,\ell}) \leq \lambda$ for all pair (j,ℓ) and and $s \in T_i$, then:

a) When $|A_i| = 0$, we have

$$\mathbb{P}(\mathcal{V}_i) \le k\lambda r + \varphi(t).$$

b) When $|A_i| > 0$, we have

$$\begin{split} &\|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}\| \\ \leq & \frac{2k^{*}\sigma}{1 - k^{2}\lambda r - \varphi(t)} \\ &+ \frac{2\sigma\sqrt{k^{2}\lambda r + \varphi(t)}(1 + (k^{*}k - k^{2})\lambda r + (k^{*} - 1)\varphi(t))}{(1 - k^{2}\lambda r - \varphi(t))^{2}} \\ &+ \frac{(2k^{*} + k)k\lambda r + (2k^{*} + 1)\varphi(t)}{1 - k^{2}\lambda r - \varphi(t)}\Delta_{\max}, \end{split}$$
where $\boldsymbol{b}_{i} := \frac{1}{|A_{i}|}\sum_{s \in A_{i}}\boldsymbol{\beta}_{s}^{*}.$

We prove Theorem 6 in Section 0c to follow. Note that Theorem 6 is similar to Theorem 4 except that the error bounds here have an additional error term $\varphi(t)$.

The procedure to derive the main Theorem 5 from Theorem 6 is exactly same as that for the Stochastic Ball Model. In particular, with r fixed to be $t\sigma\sqrt{d}$, we set $\lambda = \frac{c}{\sqrt{r\Delta_{\max}}}$. The assumptions on t and η_{\max} ensure that $\lambda k^2 r < \frac{1}{4}$ and $\varphi(t) = 2\exp(-t^2d/8) < \frac{1}{4}$. Moreover, Observations 1 and 2 still hold in the current setting. We then construct, in three steps, a partition of the fitted centers $\{\beta_i\}_{i \in [k]}$ as well as the true centers $\{\beta_s^*\}_{s \in [k^*]}$. We provide a brief description here and reuse all the notation from the proof of Theorem 3.

⁵Note that when f_s is the uniform distribution over $\mathbb{B}_s(r)$, the definition here reduces to $\rho_s(\partial_{i,j}) = \frac{1}{V_d r^d} \operatorname{ReVol}(\partial_{i,j} \cap \mathbb{B}_s(r))$ and hence is consistent with our previous definition in the ball model.

a) Step 1 (almost empty association): We consider the set

$$S_0 := \left\{ i \in [k] : \rho_s(\partial_{j,\ell}) \le \lambda, \\ \forall (s, j, \ell) \in T_i \times [k] \times [k]; |A_i| = 0 \right\}$$

Part 2(a) of Theorem 6 ensures that for all $i \in S_0$, we have $\mathbb{P}(\mathcal{V}_i) \leq k\lambda r + \varphi(t) = \frac{ck\sqrt{t}}{\sqrt{\eta_{\max}}} + \varphi(t)$. b) Step 2 (many/one fit one association):: We then

consider the fitted centers indexed by the set

$$\mathcal{J} := \{ i \in [k] : |A_i| \le 1 \} \setminus S_0.$$

For each $i \in \mathcal{J}$, there are two complementary cases:

• $\rho_s(\partial_{j,k}) \leq \lambda$ for all $(s, j, \ell) \in T_i \times [k] \times [k]$. We must have $|A_i| = 1$; say $A_i = \{s\}$. Applying Part 2(b) of Theorem 6, we have

$$\begin{split} \|\beta_{i} - \beta_{s}^{*}\| \\ = \|\beta_{i} - b_{i}\| \\ \leq \frac{2k^{*}\sigma}{1 - k^{2}\lambda r - \varphi(t)} \\ &+ \frac{2\sigma\sqrt{k^{2}\lambda r + \varphi(t)}(1 + (k^{*}k - k^{2})\lambda r + (k^{*} - 1)\varphi(t))}{(1 - k^{2}\lambda r - \varphi(t))^{2}} \\ &+ \frac{(2k^{*} + k)k\lambda r + (2k^{*} + 1)\varphi(t)}{1 - k^{2}\lambda r - \varphi(t)}\Delta_{\max} \end{split}$$

$$\stackrel{(i)}{\leq} 4k^{*}\sigma + 4\sqrt{2}\sigma(1 + c(k^{*}k - k^{2})/\sqrt{\eta_{\max}/t} + (k^{*} - 1)\varphi(t)) \\ &+ 2(c(2k^{*} + k)k/\sqrt{\eta_{\max}/t} + (2k^{*} + 1)\varphi(t))\Delta_{\max} \end{split}$$

$$\stackrel{(ii)}{\leq} \Delta_{\max} \left\{ \frac{\frac{2k^{*}}{ck\sqrt{t}} + 2ck(k^{*} + k)\sqrt{t}}{\sqrt{\eta_{\max}}} + 7k^{*}\varphi(t) \right\}$$
(49)

where in step (i) we plug in $\lambda = \frac{c}{\sqrt{\Delta_{\max}r}}$ with $r = t\sigma\sqrt{d}$ and use the assumption that $k^2\lambda r + \varphi(t) < \frac{1}{2}$; in step (ii), we use the assumption that $\eta_{\text{max}} \ge 16tc^2k^4$ to further simplify the bound.

• $\rho_s(\partial_{j,\ell}) > \lambda$ for some $(s, j, \ell) \in T_i \times [k] \times [k]$; that is, there exists some ground truth cluster \mathbb{B}_s that encloses a Voronoi boundary with a large relative volume. Applying Part 1 of Theorem 4 and plugging the value of λ , we obtain that $\|\beta_i - \beta_s^*\| \le \frac{k^*}{\lambda} + 3r \le \Delta_{\max} \frac{4k^* \sqrt{t}}{c\sqrt{\eta_{\max}}}$.

For each distinct $s \in [k^*]$ that appears in the above arguments, let $S_a^* = \{s\}$ and let the corresponding S_a index those β_i 's for which either of the two cases holds.

c) Step 3 (one-fit-many association):: We are left with the fitted centers indexed by the set

$$\mathcal{K} := \{i \in [k] : |A_i| \ge 2\} = [k] \setminus (S_0 \cup \mathcal{J}).$$

Similarly to before, for each $i \in \mathcal{K}$, there are two complementary cases:

• $\rho_s(\partial_{i,\ell}) \leq \lambda$ for all $(s,j,\ell) \in T_i \times [k] \times [k]$. Applying Part 2(b) of Theorem 4 and following the same steps as in equation (49), we obtain that

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}\| \leq \Delta_{\max} \left\{ \frac{\frac{2k^{*}}{ck\sqrt{t}} + 2ck(k^{*} + k)\sqrt{t}}{\sqrt{\eta_{\max}}} + 7k^{*}\varphi(t) \right\}$$

In this case, we let $S_a = \{i\}$ and $S_a^* = A_i$.

• $\rho_s(\partial_{i,\ell}) > \lambda$ for some $(s, j, \ell) \in T_i \times [k] \times [k]$. We then have an analogous lemma as Lemma 4 to show that β_i is close to the mean of all the true centers contained in its Voronoi set, regardless of whether we include or exclude β_s^* . The proof is exactly the same as in Section A; we omit its proof.

Lemma 12 (Proximity to mean of true centers, Gaussian case). Under the assumption of Theorem 5 let β be a local minimum of G. The following is true for each $i \in [k]$. If $\rho_s(\partial_{j,\ell}) > \lambda = \frac{c}{\sqrt{r\Delta_{\max}}} \text{ for some } (s, j, \ell) \in T_i \times [k] \times [k]$ and $|A_i \setminus \{s\}| \ge 1$, then we have the bounds

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}^{-}\| \leq \Delta_{\max} \left\{ \frac{\frac{2k^{*}}{ck\sqrt{t}} + 3ck(k^{*} + k)\sqrt{t}}{\sqrt{\eta_{\max}}} + 7k^{*}\varphi(t) \right\}$$

and

$$\|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}^{\dagger}\| \leq \Delta_{\max} \left\{ \frac{\frac{2k^{*}}{ck\sqrt{t}} + 3ck(k^{*} + k)\sqrt{t}}{\sqrt{\eta_{\max}}} + 7k^{*}\varphi(t) \right\},$$

where
$$\boldsymbol{b}_i^- := \frac{1}{|A_i \setminus \{s\}|} \sum_{s' \in A_i \setminus \{s\}} \beta_{s'}^*$$
 and $\boldsymbol{b}_i^+ := \frac{1}{|A_i \setminus \{s\}|} \sum_{s' \in A_i \setminus \{s\}} \beta_{s'}^*$

 $\overline{|A_i \cup \{s\}|} \sum_{s' \in A_i \cup \{s\}} |\mathcal{B}_{s'}^*; \text{ moreover, } |A_i \setminus \{s\}| \ge 2.$ In this case, we let $S_a = \{i\}$; also let $S_a^* = A_i \setminus \{s\}$ if the index s has appeared in the sets $\{S_{a'}^*\}$ constructed previously in Step 2 or in this step, and set $S_a^* = A_i \cup \{s\}$.

The reason that the above construction corresponds to a partition for the fitted centers $\{\beta_i\}_{i \in [k]}$ as well as a partition for the true centers $\{\beta_s^*\}_{s \in [k^*]}$ is the same as the proof in Theorem 3.

Proof of Theorem 6

To establish Theorem 6 for the Gaussian model, we follow the same strategy as in the proof of Theorem $\frac{4}{4}$ for the ball model. The only technical hurdle is that each Gaussian component distribution has unbounded support. Our main idea is to identify a bounded ball, namely $\mathbb{B}_{s}(r)$, that contains most of the probability mass of the s-th Gaussian component. Using a standard concentration inequality for χ^2 random variables (e.g., [56, Example 2.28]), we know that when t > 2, there holds the tail bound

$$\mathbb{P}_s\Big(\mathbb{B}_s(r)^{\complement}\Big) \le \varphi(t) = 2\exp(-t^2d/8),\tag{50}$$

where S^{\complement} denotes the complement of a set $S \subseteq \mathbb{R}^d$. By restricting each s-th ground truth component to the ball $\mathbb{B}_{s}(r)$ and treating the tail mass in equation (50) as additional error terms, we can repeat most of the arguments used in the proof of the ball model. In what follows, we sketch the analysis and point out the minor modifications needed to adapt the proof of Theorem $\frac{4}{4}$ to the Gaussian case.

The main step in the proof for the Ball model involves constructing smooth upper bounds for the function $W_{i \to i,s}^{\boldsymbol{v}}$ + $W_{j \to i,s}^{\boldsymbol{v}}$, as done in Proposition 3 and Proposition 4. These two propositions still hold in the Gaussian case under the i definition (48) of the "relative volume" $\rho_s(\partial_{i,j})$. In particular, the value of the integral defining $W_{i \to j,s}^{\boldsymbol{v}} + W_{j \to i,s}^{\boldsymbol{v}}$ does not increase if we restrict integration to the small set subset $(\mathbb{B}_s(r))$, as the integrand is non-positive. Consequently, we can establish the two key inequalities (18) and (19), restated below:

$$d_{i,j} \cdot \rho_s(\partial_{i,j}) \le \frac{k^*}{2} \quad \text{and} \quad \frac{D_{i,j,s}^2}{d_{i,j}} \cdot \rho_s(\partial_{i,j}) \le \frac{k^*}{2}.$$
(51)

We can then derive the structural properties of a local minimum β from the inequalities (51). As in the proof of Theorem 4, for each $i \in [k]$ indexing the fitted center β_i and its Voronoi set \mathcal{V}_i , we consider two complementary cases.

Case 1: there exist some $(s, j, \ell) \in T_i \times [k] \times [k]$ such that $\rho_s(\partial_{j,\ell}) > \lambda$.: In this case, following exactly the same argument as in the Ball model proof, we can derive from the inequalities (51) that $\|\beta_i - \beta_s^*\| \le \frac{k^*}{\lambda} + 3r$. This proves Part 1 of Theorem 6.

Case 2: for all $(s, j, \ell) \in T_i \times [k] \times [k]$ there holds $\rho_s(\partial_{j,\ell}) \leq \lambda$.: Recall that $m_{i,s}$ is the probability mass of \mathcal{V}_i with respect to the Gaussian density f_s and $c_{i,s}$ is the corresponding center of mass of \mathcal{V}_i . If we restrict the density f_s onto the ball $\mathbb{B}_s(r)$, the values of $m_{i,s}$ and $c_{i,s}$ do not change much; in particular, we can control the amount of change using the tail bound (50). With this in mind, we proceed by considering two sub cases.

Case 2(a): A_i = Ø, in which case T_i = B_i. Following the same argument for deriving equation (20) and using the Gaussian Lemma Ø in place of Lemma 8, we obtain that

$$m_{i,s} \le k\lambda r + \mathbb{P}_s \big(\mathbb{B}_s(r)^{\mathsf{C}} \big) \le k\lambda r + \varphi(t), \qquad \forall s \in B_i,$$
(52)

where the second RHS term accounts for the tail probability on $\mathbb{B}_s(r)^{\complement}$. It follows that

$$\mathbb{P}(\mathcal{V}_i) = \frac{1}{k} \sum_{s \in T_i} m_{i,s} \le k\lambda r + \varphi(t).$$

This proves Part 2(a) of Theorem 6.

• Case 2(b): $A_i \neq \emptyset$. By Lemma 2, β must satisfy the expression

$$\boldsymbol{\beta}_i = \frac{\sum_{s \in [k^*]} \boldsymbol{c}_{i,s} m_{i,s}}{\sum_{s \in [k^*]} m_{i,s}}$$

Using this expression, we have the decomposition $\beta_i - b_i = \mu + \nu$ for some vectors μ and ν as in equation (22). To bound μ and ν , we follow our general strategy to decompose the Gaussian density f_s into two parts, one supported on the ball $\mathbb{B}_s(r)$ and the other the tail, where the tail probability is bounded by $\varphi(t)$ as in equation (50). By doing so, we can establish analogous versions of the bounds (29) and (30) as given below:

and

$$\|\boldsymbol{\beta}_{s}^{*}-\boldsymbol{c}_{i,s}\| \leq \begin{cases} \frac{2\sigma}{m_{i,s}}, & \text{if } s \in B_{i}, \\ \frac{2\sigma\sqrt{k^{2}\lambda r+\varphi(t)}}{m_{i,s}}, & \text{if } s \in A_{i}, \end{cases}$$

 $\begin{cases} m_{i,s} \le k\lambda r + \varphi(t), & \text{if } s \in B_i, \\ m_{i,s} \ge 1 - k^2\lambda r - \varphi(t), & \text{if } s \in A_i, \end{cases}$

where the bound on $\|\beta_s^* - c_{i,s}\|$ follows from Lemma [1]. Using the above two bounds, we can further establish analogous versions of Lemmas [5] and [6] as given below:

$$\|\boldsymbol{\mu}\| \leq \frac{2k^*\sigma}{1 - k^2\lambda r - \varphi(t)}$$

$$+ \frac{2k^*(k\lambda r + \varphi(t))\sigma\sqrt{k^2\lambda r + \varphi(t)}}{(1 - k^2\lambda r - \varphi(t)))^2} + \frac{2k^*(k\lambda r + \varphi(t))}{1 - k^2\lambda r - \varphi(t)}\Delta_{\max}, \|\boldsymbol{\nu}\| \le \frac{2\sigma\sqrt{k^2\lambda r + \varphi(t)}}{1 - k^2\lambda r - \varphi(t)} + \frac{k^2\lambda r + \varphi(t)}{1 - k^2\lambda r - \varphi(t)}\Delta_{\max}.$$

It follows that an analogue of inequality (??) holds:

$$\begin{split} \|\boldsymbol{\beta}_{i} - \boldsymbol{b}_{i}\| \\ \leq \|\boldsymbol{\mu}\| + \|\boldsymbol{\nu}\| \\ \leq \frac{2k^{*}\sigma}{1 - k^{2}\lambda r - \varphi(t)} \\ + \frac{2\sigma\sqrt{k^{2}\lambda r + \varphi(t)}(1 + (k^{*}k - k^{2})\lambda r + (k^{*} - 1)\varphi(t))}{(1 - k^{2}\lambda r - \varphi(t))^{2}} \\ + \frac{(2k^{*} + k)k\lambda r + (2k^{*} + 1)\varphi(t)}{1 - k^{2}\lambda r - \varphi(t)}\Delta_{\max}. \end{split}$$

This proves Part 2(b) of Theorem 6.

ADDITIONAL EXAMPLES

In this section, we provide two concrete examples that give additional insights on the behaviors of the local minima of the k-means objective and corroborate the results in our main theorems.

Our main theorem assumes certain separation conditions in terms of the SNRs η_{\min} and η_{\max} . The first example shows that if the SNR is too small, then a local minimum may fail to have the structure described in Theorem [3]. Therefore, a separation condition on the true clusters is in general necessary.

Example 1 (Small Separation). Consider the Stochastic Ball Model with k = 3 in dimension d = 1, where the ground truth cluster centers are $\beta_1^* = -1$, $\beta_2^* = 0$ and $\beta_3^* = 1$, and the radius r of the balls satisfies $(\frac{9\sqrt{2}}{2} - \frac{1}{4})r > 1$. Let $\beta = (\beta_1, \beta_2, \beta_3)$ be a candidate solution with $\beta_1 = -\frac{2}{3} - \frac{1}{6}r$, $\beta_2 = \frac{2}{3} + \frac{1}{6}r$ and $\beta_3 > 0$ sufficiently large.

When β_3 is sufficiently large, the minimization $\min_{i \in [k]} ||x - \beta_i||^2$ in the objective *G* is never attained by i = 3. In this case, the only effective variables for *G* are the first two centers β_1 and β_2 . The Voronoi boundary $\partial_{1,2}(\beta)$ (which is 0) intersects the second ground truth cluster. Note that these properties continue to hold under small perturbation of β . Consequently, for any solution *b* in a small neighborhood of β , its objective value has the following expression:

$$G(\mathbf{b}) = \frac{1}{6r} \left[\int_{-1-r}^{-1+r} (x-b_1)^2 dx + \int_{-r}^{\frac{b_1+b_2}{2}} (x-b_1)^2 dx + \int_{\frac{b_1+b_2}{2}}^{r} (x-b_2)^2 dx + \int_{1-r}^{1+r} (x-b_2)^2 dx \right].$$

We compute the gradient and Hessian of G at b (recall that only the first two coordinate of b are effective):

$$\nabla_{\boldsymbol{b}}G = \frac{1}{6r} \left[\begin{array}{c} -2\int_{-1-r}^{-1+r} (x-b_1) \mathrm{d}x - 2\int_{-r}^{\frac{b_1+b_2}{2}} (x-b_1) \mathrm{d}x \\ -2\int_{\frac{b_1+b_2}{2}}^{r} (x-b_2) \mathrm{d}x - 2\int_{1-r}^{1+r} (x-b_2) \mathrm{d}x \end{array} \right],$$



Figure 7. Two-dimensional Stochastic Ball Model, where the true centers are $\beta_1^* = (-1,0)$, $\beta_2^* = (0,0)$ and $\beta_3^* = (1,0)$. The third fitted center β_3 is far away from the origin, so the only effective variables are β_1 and β_2 . Left panel: $\beta_1 = \beta_1^*$ and $\beta_2 = \frac{1}{2}(\beta_2^* + \beta_3^*)$. The green line is the Voronoi boundary $\partial_{1,2}(\beta)$. Right panel: A perturbed solution (β_1', β_2') and the corresponding Voronoi boundary.

$$\nabla_b^2 G = \frac{1}{6r} \begin{bmatrix} 6r + (b_1 + b_2) - \frac{b_2 - b_1}{2} & -\frac{b_2 - b_1}{2} \\ -\frac{b_2 - b_1}{2} & 6r - (b_1 + b_2) + \frac{b_1 - b_2}{2} \end{bmatrix}$$

Evaluating these expressions at $\beta_1 = -\frac{2}{3} - \frac{1}{6}r$ and $\beta_2 = \frac{2}{3} + \frac{1}{6}r$, we find that the gradient vanishes $\nabla_{\boldsymbol{b}}G \mid_{\boldsymbol{b}=\boldsymbol{\beta}} = 0$ and the Hessian is

$$\nabla_{\boldsymbol{b}}^2 G \mid_{\boldsymbol{b}=\boldsymbol{\beta}} = \frac{1}{6r} \begin{bmatrix} \frac{35}{6}r - \frac{2}{3} & -\frac{2}{3} - \frac{1}{6}r \\ -\frac{2}{3} - \frac{1}{6}r & \frac{37}{6}r + \frac{2}{3} \end{bmatrix}.$$

When $(\frac{9\sqrt{2}}{2} - \frac{1}{4})r > 1$ or equivalently $\eta_{\min} < \frac{9\sqrt{2}}{2} - \frac{1}{4}$, the Hessian is positive definite, so β is a local minimum of G. Moreover, one can verify that $G(\beta) < G(\beta^*)$, so β is not a global minimum. We see that the spurious local minimum β does not have the structure described in Theorem 3 as β involves a 2-fit-3 association.

The second example shows that in higher dimensions, there exists a local minimum $\beta = (\beta_1, \beta_2, \beta_3)$ such that β_1 approximately equals β_1^* , and β_2 approximately equals $(\beta_2^* + \beta_3^*)/2$ — a structure guaranteed by Theorem 3 — but neither approximation is exact. Therefore, the non-zero approximation errors that appear in Theorem 3, is necessary in general.

Example 2 (Approximation Errors). Consider the Stochastic Ball model with k = 3 in dimension d = 2, where the true cluster centers are $\beta_1^* = (-1,0)$, $\beta_2^* = (0,0)$ and $\beta_3^* = (1,0)$, where the radius r of the balls satisfies $r \ge \frac{1}{4}$. Let β be a candidate solution with $\beta_1 = (-1,0)$, $\beta_2 = (\frac{1}{2},0)$ and β_3 sufficiently far away from the origin. See the left panel of Figure 7 for an illustration.

As in Example 1, here the only effective variables are β_1 and β_2 . Assume first that $r = \frac{1}{4}$. In this case, the Voronoi boundary $\partial_{1,2}$ is at $x_1 = -\frac{1}{4} = \dot{\beta}_2^* - r$, the left boundary of the second true cluster. We claim that β is a local minimum of G; the proof is deferred to the end of this section. Now, let us increase the radius r by a sufficient small amount, in which case the objective function becomes G. By the continuity, there exists a local minimum $\tilde{\beta}$ of \tilde{G} near the original local minima β . Recall that by Lemma 2, β_1 and β_2 must lie at the center of mass of their Voronoi sets $\mathcal{V}_1(\boldsymbol{\beta})$ and $\mathcal{V}_2(\boldsymbol{\beta})$, respectively. It is then not hard to see that the new Voronoi boundary $\partial_{1,2}$ corresponding to $\hat{\beta}$ necessarily intersects the interior of the second true cluster \mathbb{B}_2 . It follows that $\mathcal{V}_1(\boldsymbol{\beta}) = \mathbb{B}_1 \cup D$ and $\mathcal{V}_2(\beta) = (\mathbb{B}_2 \cup \mathbb{B}_3) \setminus D$ for some subset $D \subset \mathbb{B}_2$ with a positive measure. Applying Lemma 2 again, we conclude that β_1 is close but not equal to β_1^* , and that β_2 is close but not equal to $(\beta_1^* + \beta_2^*)/2$.

Proof of the claim. Let $t \in (0, 1/8)$ be a sufficiently small number, and $v_1, v_2 \in \mathbb{R}^2$ be two arbitrary vectors satisfying $||v_1||^2 + ||v_2||^2 = 1$. Consider perturbing β_1 and β_2 to $\beta'_1 = \beta_1 + tv_1$ and $\beta'_2 = \beta_2 + tv_2$, respectively. Since Voronoi sets only change by a small amount when the perturbation t is small, we find that $\Delta^v_{2\to 1}(t) \subseteq \mathbb{B}_2$ is the only set of points that change their association from one Voronoi set to another; see the right panel of Figure 7. Using the expression (10) for the directional k-means objective, we can write $G(\beta')$ as

$$\begin{split} G(\boldsymbol{\beta}') = & H^{\boldsymbol{v}}(t) \\ = & \frac{1}{3} \bigg[\int_{\mathcal{V}_1(\boldsymbol{\beta}) \cap \mathbb{B}_1} \|\boldsymbol{x} - \boldsymbol{\beta}_1 - t\boldsymbol{v}_1\|^2 \mathrm{d}\boldsymbol{x} \\ & + \int_{\mathcal{V}_2(\boldsymbol{\beta}) \cap (\mathbb{B}_2 \cup \mathbb{B}_3)} \|\boldsymbol{x} - \boldsymbol{\beta}_2 - t\boldsymbol{v}_2\|^2 \mathrm{d}\boldsymbol{x} \bigg] \\ & + \frac{1}{3} \int_{\Delta_{2 \to 1}^{\boldsymbol{v}}(t)} \left(\|\boldsymbol{x} - \boldsymbol{\beta}_1 - t\boldsymbol{v}_1\|^2 - \|\boldsymbol{x} - \boldsymbol{\beta}_2 - t\boldsymbol{v}_2\|^2 \right) \mathrm{d}\boldsymbol{x}. \end{split}$$

A quick calculation shows that

$$G(\boldsymbol{\beta}') - G(\boldsymbol{\beta}) = H^{\boldsymbol{v}}(t) - H^{\boldsymbol{v}}(0) = \frac{1}{3} \bigg[t^2 - \underbrace{\int_{\Delta_{2\to1}^{\boldsymbol{v}}(t)} \left(\|\boldsymbol{x} - \boldsymbol{\beta}_2 - t\boldsymbol{v}_1\|^2 - \|\boldsymbol{x} - \boldsymbol{\beta}_1 - t\boldsymbol{v}_2\|^2 \right) \mathrm{d}\boldsymbol{x}}_{K(t)} \bigg].$$

We decompose the term K(t) as follows:

$$\begin{split} K(t) = & \int_{\Delta_{2 \to 1}^{\boldsymbol{v}}(t)} \underbrace{\left(\|\boldsymbol{x} - \boldsymbol{\beta}_{2}\|^{2} - \|\boldsymbol{x} - \boldsymbol{\beta}_{1}\|^{2} \right)}_{\kappa_{1}} \mathrm{d}\boldsymbol{x} \\ &+ \int_{\Delta_{2 \to 1}^{\boldsymbol{v}}(t)} \underbrace{t^{2} \left(\|\boldsymbol{v}_{2}\|^{2} - \|\boldsymbol{v}_{1}\|^{2} \right)}_{\kappa_{2}} \mathrm{d}\boldsymbol{x} \\ &+ \int_{\Delta_{2 \to 1}^{\boldsymbol{v}}(t)} \underbrace{2t \left(\langle \boldsymbol{v}_{1} - \boldsymbol{v}_{2}, \boldsymbol{x} \rangle + \langle \boldsymbol{v}_{2}, \boldsymbol{\beta}_{2} \rangle - \langle \boldsymbol{v}_{1}, \boldsymbol{\beta}_{1} \rangle \right)}_{\kappa_{3}} \mathrm{d}\boldsymbol{x} \end{split}$$

For all $\boldsymbol{x} \in \Delta_{2 \to 1}^{\boldsymbol{v}}(t) \subseteq \mathcal{V}_2(\boldsymbol{\beta})$, we have $\|\boldsymbol{x} - \boldsymbol{\beta}_2\| \leq \|\boldsymbol{x} - \boldsymbol{\beta}_1\|$, hence $\kappa_1 \leq 0$. We also have $\kappa_2 \leq t^2$ since $\|\boldsymbol{v}_2\|^2 - \|\boldsymbol{v}_1\|^2 \leq \|\boldsymbol{v}_1\|^2 \leq \|\boldsymbol{v}_1\|^2 \leq 1$. To bound κ_3 we observe that $\langle \boldsymbol{v}_1 - \boldsymbol{v}_2, \boldsymbol{x} \rangle \leq \|\boldsymbol{v}_1 - \boldsymbol{v}_2\| \|\boldsymbol{x}\| \leq 4r = 1$ for all $\boldsymbol{x} \in \Delta_{2 \to 1}^{\boldsymbol{v}}(t) \subseteq \mathbb{B}_2, \langle \boldsymbol{\beta}_2, \boldsymbol{v}_2 \rangle \leq \|\boldsymbol{\beta}_2\| \|\boldsymbol{v}_2\| \leq \frac{1}{2}$, and $\langle \boldsymbol{\beta}_1, \boldsymbol{v}_1 \rangle \leq \|\boldsymbol{\beta}_1\| \|\boldsymbol{v}_1\| \leq 1$; it follows that $\kappa_3 \leq 5t$. Combining pieces, we obtain that

$$G(\boldsymbol{\beta}') - G(\boldsymbol{\beta}) \ge \left[t^2 - 6t \cdot \operatorname{Vol}(\Delta_{2 \to 1}^{\boldsymbol{v}}(t))\right]/3.$$
(53)

It remains to control the volume of the set $\Delta_{2\to1}^{v}(t)$, which is illustrated in Figure 8 Observe that the distance between the old mid point $(\beta_1 + \beta_2)/2$ and the new Voronoi boundary $\partial_{1,2}(\beta')$ can be bounded as

$$d_{1} := \operatorname{dist}\left(\frac{\beta_{1} + \beta_{2}}{2}, \partial_{1,2}(\beta')\right)$$
$$\leq \left\|\frac{\beta_{1} + \beta_{2}}{2} - \frac{\beta_{1}' + \beta_{2}'}{2}\right\|$$
$$\leq t.$$

Moreover, the (unsigned) angle ψ between the old and new Voronoi boundaries $\partial_{1,2}(\beta)$ and $\partial_{1,2}(\beta')$ satisfies

$$\tan \psi = \left| \frac{t(v_{2,2} - v_{1,2})}{2 + t(v_{2,1} - v_{1,1})} \right| \le \frac{t}{1 - t} \le 2t.$$



Figure 8. Illustration of $\Delta_{2\to1}^{v}(t)$. The new boundary $\partial_{1,2}(\beta')$ has an angle ψ with the original Voronoi boundary $\partial_{1,2}(\beta)$.

From these two observations and the fact that r = 1/4, elementary geometry shows that the distance d_2 between β_2^* and $\partial_{1,2}(\beta')$ satisfies

$$d_{2} = r \cos \psi - d_{1}$$

= $\frac{r}{\sqrt{1 + \tan^{2} \psi}} - d_{1} \ge \frac{r}{1 + 2t} - 4r \ge r(1 - 6t),$

whence

$$\operatorname{Vol}(\Delta_{2 \to 1}^{v}(t)) \le 2 \cdot \sqrt{r^2 - d_2^2} \cdot (r - d_2) \le 12t\sqrt{12t}$$

Combining with equation (53) shows that $G(\beta') > G(\beta)$ when t is sufficiently small. As this inequality holds for arbitrary perturbation direction (v_1, v_2) , we conclude that β is a local minimum of G.

REFERENCES

- A. K. Jain, "Data clustering: 50 years beyond k-means," *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651–666, 2010.
- [2] S. Dasgupta, "The hardness of k-means clustering," Department of Computer Science and Engineering, University of California, San Diego, Tech. Rep., 2008.
- [3] M. Mahajan, P. Nimbhorkar, and K. Varadarajan, "The planar k-means problem is NP-hard," in *International Workshop on Algorithms and Computation*. Springer, 2009, pp. 274–285.
- [4] P. Awasthi, M. Charikar, R. Krishnaswamy, and A. K. Sinop, "The hardness of approximation of Euclidean kmeans," in 31st International Symposium on Computational Geometry. arXiv preprint arXiv:1502.03316, 2015, pp. 754–767.
- [5] J. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, no. 14. Oakland, CA, USA, 1967, pp. 281–297.
- [6] J. A. Hartigan and M. A. Wong, "Algorithm AS 136: A k-means clustering algorithm," *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, vol. 28, no. 1, pp. 100–108, 1979.
- [7] D. Steinley, "Local optima in k-means clustering: what you don't know may hurt you." *Psychological Methods*, vol. 8, no. 3, p. 294, 2003.

- [8] —, "Profiling local optima in k-means clustering: Developing a diagnostic technique." *Psychological methods*, vol. 11, no. 2, p. 178, 2006.
- [9] K. Chaudhuri, S. Dasgupta, and A. Vattani, "Learning mixtures of Gaussians using the k-means algorithm," arXiv preprint arXiv:0912.0086, 2009.
- [10] J. Xu, D. J. Hsu, and A. Maleki, "Global analysis of expectation maximization for mixtures of two Gaussians," in Advances in Neural Information Processing Systems, 2016, pp. 2676–2684.
- [11] C. Daskalakis, C. Tzamos, and M. Zampetakis, "Ten steps of em suffice for mixtures of two gaussians," in *Conference on Learning Theory*. PMLR, 2017, pp. 704– 710.
- [12] C. Jin, Y. Zhang, S. Balakrishnan, M. J. Wainwright, and M. I. Jordan, "Local maxima in the likelihood of Gaussian mixture models: Structural results and algorithmic consequences," in *Advances in Neural Information Processing Systems*, 2016, pp. 4116–4124.
- [13] S. Lloyd, "Least squares quantization in PCM," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, 1982.
- [14] L. Bottou and Y. Bengio, "Convergence properties of the k-means algorithms," in *Advances in Neural Information Processing Systems*, 1995, pp. 585–592.
- [15] N. Srebro, "Are there local maxima in the infinitesample likelihood of Gaussian mixture estimation?" in *International Conference on Computational Learning Theory*, 2007, pp. 628–629.
- [16] D. Steinley, "K-means clustering: a half-century synthesis," *British Journal of Mathematical and Statistical Psychol*ogy, vol. 59, no. 1, pp. 1–34, 2006.
- [17] A. Kumar, Y. Sabharwal, and S. Sen, "A simple linear time $(1+\epsilon)$ -approximation algorithm for k-means clustering in any dimensions," in *Annual Symposium on Foundations of Computer Science*, vol. 45. IEEE Computer Society Press, 2004, pp. 454–462.
- [18] T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "A local search approximation algorithm for k-means clustering," *Computational Geometry*, vol. 28, no. 2-3, pp. 89–112, 2004.
- [19] G. W. Milligan, "An examination of the effect of six types of error perturbation on fifteen clustering algorithms," *Psychometrika*, vol. 45, no. 3, pp. 325–342, 1980.
- [20] S. Har-Peled and B. Sadri, "How fast is the k-means method?" Algorithmica, vol. 41, no. 3, pp. 185–202, 2005.
- [21] D. Arthur and S. Vassilvitskii, "How slow is the k-means method?" in *Symposium on Computational Geometry*, vol. 6, no. 32, 2006, pp. 1–10.
- [22] A. Kumar and R. Kannan, "Clustering with spectral norm and the k-means algorithm," in 2010 IEEE 51st Annual Symposium on Foundations of Computer Science. IEEE, 2010, pp. 299–308.
- [23] Y. Lu and H. H. Zhou, "Statistical and computational guarantees of Lloyd's algorithm and its variants," arXiv preprint arXiv:1612.02099, 2016.
- [24] D. Arthur and S. Vassilvitskii, "k-means++: The advantages of careful seeding," in *Proceedings of the Eighteenth*

Annual ACM-SIAM Symposium on Discrete algorithms. Society for Industrial and Applied Mathematics, 2007, pp. 1027–1035.

- [25] R. Ostrovsky, Y. Rabani, L. J. Schulman, and C. Swamy, "The effectiveness of Lloyd-type methods for the k-means problem," *Journal of the ACM (JACM)*, vol. 59, no. 6, p. 28, 2012.
- [26] S. Dasgupta and L. Schulman, "A probabilistic analysis of EM for mixtures of separated, spherical Gaussians," *Journal of Machine Learning Research*, vol. 8, no. Feb, pp. 203–226, 2007.
- [27] J. Peng and Y. Xia, "A new theoretical framework for k-means-type clustering," in *Foundations and Advances in Data Mining*. Springer, 2005, pp. 79–96.
- [28] J. Peng and Y. Wei, "Approximating k-means-type clustering via semidefinite programming," *SIAM Journal on Optimization*, vol. 18, no. 1, pp. 186–205, 2007.
- [29] E. Elhamifar, G. Sapiro, and R. Vidal, "Finding exemplars from pairwise dissimilarities via simultaneous sparse recovery," in *Advances in Neural Information Processing Systems*, 2012, pp. 19–27.
- [30] A. Nellore and R. Ward, "Recovery guarantees for exemplar-based clustering," *Information and Computation*, vol. 245, pp. 165–180, 2015.
- [31] P. Awasthi, A. S. Bandeira, M. Charikar, R. Krishnaswamy, S. Villar, and R. Ward, "Relax, no need to round: Integrality of clustering formulations," in *Proceedings* of the 2015 Conference on Innovations in Theoretical Computer Science. ACM, 2015, pp. 191–200.
- [32] X. Li, Y. Li, S. Ling, T. Strohmer, and K. Wei, "When do birds of a feather flock together? k-means, proximity, and conic programming," *Mathematical Programming*, vol. 179, no. 1-2, pp. 295–341, 2020.
- [33] Y. Fei and Y. Chen, "Hidden integrality of SDP relaxation for sub-Gaussian mixture models," in *Conference on Learning Theory (COLT)*, 2018, arXiv preprint arXiv:1803.06510.
- [34] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 39, no. 1, pp. 1–22, 1977.
- [35] S. Balakrishnan, M. J. Wainwright, and B. Yu, "Statistical guarantees for the EM algorithm: From population to sample-based analysis," *The Annals of Statistics*, vol. 45, no. 1, pp. 77–120, 2017.
- [36] W. Qian, Y. Zhang, and Y. Chen, "Global convergence of least squares EM for demixing two log-concave densities," in Advances in Neural Information Processing Systems, 2019, pp. 4795–4803.
- [37] J. Kwon, W. Qian, C. Caramanis, Y. Chen, and D. Davis, "Global convergence of the EM algorithm for mixtures of two component linear regression," in *Conference on Learning Theory*, 2019, pp. 2055–2110.
- [38] J. Xu, D. J. Hsu, and A. Maleki, "Benefits of overparameterization with EM," in *Advances in Neural Information Processing Systems*, 2018, pp. 10662–10672.
- [39] X. Yi and C. Caramanis, "Regularized EM algorithms: A unified framework and statistical guarantees," in Advances

in Neural Information Processing Systems, 2015, pp. 1567–1575.

- [40] J. M. Klusowski, D. Yang, and W. D. Brinda, "Estimating the coefficients of a mixture of two linear regressions by expectation maximization," *IEEE Transactions on Information Theory*, 2019.
- [41] B. Yan, M. Yin, and P. Sarkar, "Convergence of gradient EM on multi-component mixture of Gaussians," in Advances in Neural Information Processing Systems, 2017, pp. 6956–6966.
- [42] J. Kwon and C. Caramanis, "EM converges for a mixture of many linear regressions," arXiv preprint arXiv:1905.12106, 2019.
- [43] R. Kannan, H. Salmasian, and S. Vempala, "The spectral method for general mixture models," in *International Conference on Computational Learning Theory*. Springer, 2005, pp. 444–457.
- [44] S. Mei, Y. Bai, and A. Montanari, "The landscape of empirical risk for non-convex losses," *The Annals of Statistics*, vol. 46, no. 6A, pp. 2747–2774, 2018.
- [45] S. Bhojanapalli, B. Neyshabur, and N. Srebro, "Global optimality of local search for low rank matrix recovery," in Advances in Neural Information Processing Systems, 2016, pp. 3873–3881.
- [46] I. Safran and O. Shamir, "Spurious local minima are common in two-layer ReLU neural networks," in *International Conference on Machine Learning*, 2018, pp. 4430–4438.
- [47] R. Ge, J. D. Lee, and T. Ma, "Matrix completion has no spurious local minimum," in Advances in Neural Information Processing Systems 29, 2016.
- [48] G. H. Ball and D. J. Hall, "PROMENADE an on-line pattern recognition system," Stanford Research Institution, Menlo Park, CA, Tech. Rep., 1967.
- [49] M. M. Astrahan, "Speech analysis by clustering, or the hyperphoneme method," Stanford University, Department of Computer Science, Tech. Rep., 1970.
- [50] A. R. Barakbah and A. Helen, "Optimized k-means: an algorithm of initial centroids optimization for k-means," in *Proc. Seminar on Soft Computing, Intelligent System,* and Information Technology (SIIT), Surabaya, 2005.
- [51] A. R. Barakbah and Y. Kiyoki, "A pillar algorithm for k-means optimization by distance maximization for initial centroid designation," in 2009 IEEE Symposium on Computational Intelligence and Data Mining. IEEE, 2009, pp. 61–68.
- [52] P. Fränti and S. Sieranoja, "K-means properties on six clustering benchmark datasets," *Applied Intelligence*, vol. 48, no. 12, pp. 4743–4759, 2018.
- [53] R.-D. Buhai, A. Risteski, Y. Halpern, and D. Sontag, "Benefits of overparameterization in single-layer latent variable generative models," *arXiv preprint arXiv:1907.00030*, 2019.
- [54] R. Dwivedi, N. Ho, K. Khamaru, M. I. Jordan, M. J. Wainwright, and B. Yu, "Singularity, misspecification, and the convergence rate of EM," *arXiv preprint arXiv:1810.00828*, 2018.
- [55] Y. Zhang, H.-W. Kuo, and J. Wright, "Structured local minima in sparse blind deconvolution," in *Advances in*

Neural Information Processing Systems 31, 2018, pp. 2328–2337.

[56] M. J. Wainwright, *High-Dimensional Statistics: A Non-asymptotic Viewpoint*. Cambridge University Press, 2019, vol. 48.

Wei Qian received her Ph.D. and M.S. degrees in Operations Research from Cornell University in 2020, and B.S. degree in mathematics from The University of Michigan, Ann Arbor in 2014. Her research work lie in machine learning, reinforcement learning and optimization, with applications in transportation systems.

Yuqian Zhang is an Assistant Professor with the Department of Electrical and Computer Engineering at Rutgers University. She was a postdoctoral scholar with the Tripods Center for Data Science at Cornell University. She obtained her Ph.D. and M.S. in Electrical Engineering from Columbia University, and B.S. in Information Engineering from Xi'an Jiaotong University. Her research leverages physical models in data driven computations, convex and nonconvex optimization, solving problems in machine learning, computer vision, signal processing.

Yudong Chen is an Associate Professor with the School of Operations Research and Information Engineering at Cornell University. He obtained his Ph.D. degree in Electrical and Computer Engineering in 2013 from The University of Texas at Austin, and M.S. and B.S. degrees in Control Science and Engineering from Tsinghua University. He was a postdoctoral scholar in the Electrical Engineering and Computer Sciences department at the University of California, Berkeley from 2013 to 2015. He has served as area chairs for AAAI, AISTATS and NeurIPS, and received a National Science Foundation CAREER award. His research work lies in machine learning, reinforcement learning, high-dimensional statistics, and optimization, with applications in network scheduling, wireless communication, and financial systems.