

**Organization.** In this supplementary material, we include the technical sections in its entirety (and not just the proofs of the lemma and theorem statements in the main body), for convenience. First we provide a more comprehensive literature review and review of prior work on mixtures of Gaussians. In Section 7, we introduce the notation in more detail, and provide proofs of various properties that will be useful in both the upper bound and lower bound. In Section 8, we prove that the  $k$ -means++ algorithm (or Lloyd’s with appropriate initialization) recovers most of the cluster memberships correctly. In Section 9, we show that any locally optimal algorithm (including Lloyd’s algorithm) makes at least  $\Omega(kd/\Delta^4)$  points.

## 6. Introduction and Related Work

**Considerations in the choice of the Semi-random GMM model.** Here we briefly discuss different semi-random models, and considerations involved in favoring Definition 1.1. Another semi-random model that comes to mind is one that can move each point closer to the mean of its own cluster (closer just in terms of distance, regardless of direction). Intuitively this seems appealing since this improves the cost of the planted clustering. However, in this model the optimal  $k$ -means clustering of the perturbed instance can be vastly different from the planted solution. This is because one can move many points  $x$  in cluster  $C_i$  in such a way that  $x$  becomes closer to a different mean rather than  $\mu_i$ . For high dimensional Gaussians it is easy to see that the distance of each point to its own mean will be on the order of  $(\sqrt{d} + 2\sqrt{\log N})\sigma$ . Hence, in our regime of interest, the inter mean separation of  $\sqrt{k \log N}\sigma$  could be much smaller than the radius of any cluster (when  $d \gg k$ ). Consider an adversary that moves a large fraction of the points in a given cluster to the mean of another cluster. While the distance of these points to their cluster mean has only decreased from roughly  $(\sqrt{d} + 2\sqrt{\log N})\sigma$  to around  $\sqrt{k \log N}\sigma$ , these points now become closer to the mean of a different cluster! In the semi-random GMM model on the other hand, the adversary is only allowed to move the point  $x$  along the direction of  $x - \mu_i$ ; hence, each point  $x$  becomes closer to its own mean than to the means of other clusters. Our results show that in such a model, the optimal clustering solution can change by at most  $\tilde{O}(d/\Delta^4)$  points.

### 6.1. Related Work

There has been a long line of algorithmic results on Gaussian mixture models starting from (Teicher, 1961; 1967; Pearson, 1894). These results fall into two broad

categories: (1) *Clustering algorithms*, which aim to recover the component/cluster memberships of the points and (2) *Parameter estimation*, where the goal is to estimate the parameters of the Gaussian components. When the components of the mixture are sufficiently well-separated, i.e.,  $\|\mu_i - \mu_j\|_2 \geq \sigma\sqrt{\log(Nk)}$ , then the Gaussians do not overlap w.h.p., and then the two tasks become equivalent w.h.p. We now review the different algorithms that have been designed for these two tasks, and comment on their robustness to semi-random perturbations.

**Clustering Algorithms.** The first polynomial time algorithmic guarantees were given by Dasgupta (Dasgupta, 1999), who showed how to cluster a mixture of  $k$  Gaussians with identical covariance matrices when the separation between the cluster means is of the order  $\Omega(\sigma\sqrt{d}\text{polylog}(N))$ , where  $\sigma$  denotes the maximum variance of any cluster along any direction<sup>4</sup>. Distance-based clustering algorithms that are based on strong distance-concentration properties of high-dimensional Gaussians improved the separation requirement between means  $\mu_i$  and  $\mu_j$  to be  $\Omega(d^{1/4}\text{polylog}(N))(\sigma_i + \sigma_j)$  (Arora & Kannan, 2001; Dasgupta & Schulman, 2007), where  $\sigma_i$  denotes the maximum variance of points in cluster  $i$  along any direction. Vempala and Wang (Vempala & Wang, 2004) and subsequent results (Kannan et al., 2008; Achlioptas & McSherry, 2005) used PCA to project down to  $k$  dimensions (when  $k \leq d$ ), and then used the above distance-based algorithms to get state-of-the-art guarantees for many settings: for spherical Gaussians a separation of roughly  $\|\mu_i - \mu_j\|_2 \geq (\sigma_i + \sigma_j) \min\{k, d\}^{1/4} \text{polylog}(N)$  suffices (Vempala & Wang, 2004). For non-spherical Gaussians, a separation of  $\|\mu_i - \mu_j\|_2 \geq (\sigma_i + \sigma_j)k^{3/2}\sqrt{\log N}$  is known to suffice (Achlioptas & McSherry, 2005; Kannan et al., 2008). Brubaker and Vempala (Brubaker & Vempala, 2008) gave a qualitative improvement on the separation requirement for non-spherical Gaussians by having a dependence only on the variance along the direction of the line joining the respective means, as opposed to the maximum variance along any direction.

Recent work has also focused on provable guarantees for heuristics such as the Lloyd’s algorithm for clustering mixtures of Gaussians (Kumar & Kannan, 2010; Awasthi & Sheffet, 2012). Iterative algorithms like the Lloyd’s algorithm (also called  $k$ -means algorithm) (Lloyd, 1982) and its variants like  $k$ -means++ (Ostrovsky et al., 2006; Arthur & Vassilvitskii, 2007) are the method-of-choice for clustering in practice. The best known guarantee (Awasthi &

<sup>4</sup>The  $\text{polylog}(N)$  term involves a dependence of either  $(\log N)^{1/4}$  or  $(\log N)^{1/2}$ .

Sheffet, 2012) along these lines requires a separation of order  $\sigma\sqrt{k\log N}$  between any pair of means, where  $\sigma$  is the maximum variance among all clusters along any direction. To summarize, for a mixture of  $k$  Gaussians in  $d$  dimensions with variance of each cluster being bounded by  $\sigma^2$  in every direction, the state-of-the-art guarantees require a separation of roughly  $\sigma \min\{k, d\}^{1/4} \text{polylog}(N)$  between the means of any two components (Vempala & Wang, 2004) for spherical Gaussians, while a separation of  $\sigma\sqrt{\min\{k, d\}\log N}$  is known to suffice for non-spherical Gaussians (Awasthi & Sheffet, 2012).

The techniques in many of the above works rely on strong distance concentration properties of high-dimensional Gaussians. For instance, the arguments of (Arora & Kannan, 2001; Vempala & Wang, 2004) that obtain a separation of order  $\min\{k^{1/4}, d^{1/4}\}$  crucially rely on the tight concentration of the squared distance around  $\sigma^2(d \pm c\sqrt{d})$ , between any pair of points in the same cluster. These arguments do not seem to carry over to the semi-random model. Brubaker (Brubaker, 2009) gave a robust algorithm for clustering a mixture of Gaussians when at most  $o(1/k)$  fraction of the points are corrupted arbitrarily. However, it is unclear if the arguments can be modified to work under the semi-random model, since the perturbations can potentially affect all the points in the instance. On the other hand, our results show that the Lloyd’s algorithm of Kumar and Kannan (Kumar & Kannan, 2010) is robust to these semi-random perturbations.

Finally, there has also been significant work on designing clustering algorithms under deterministic assumptions on the data such as resilience of the optimal clustering to distance perturbations (Ackerman & Ben-David, 2009; Awasthi et al., 2012; Balcan & Liang, 2012; Angelidakis et al., 2017; Dutta et al., 2017). The assumptions in these works are incomparable to those in our work and, in particular the separation requirement is more stringent when applied to the special case of the Gaussian mixture models.

**Parameter Estimation.** A different approach is to design algorithms that estimate the parameters of the underlying Gaussian mixture model, and then assuming the means are well separated, accurate clustering can be performed. A very influential line of work focuses on the method-of-moments (Kalai et al., 2010; Moitra & Valiant, 2010; Belkin & Sinha, 2010) to learn the parameters of the model when the number of clusters  $k = O(1)$ . Moment methods (necessarily) require running time (and sample complexity) of roughly  $d^{O(k^2)}$ , but do not assume any explicit separation between the components of the mixture. Recent work (Hsu

& Kakade, 2013; Bhaskara et al., 2014b; Goyal et al., 2014; Bhaskara et al., 2014a; Anderson et al., 2014; Ge et al., 2015) uses uniqueness of tensor decompositions (of order 3 and above) to implement the method of moments and give polynomial time algorithms assuming the means are sufficiently high dimensional, and do not lie in certain degenerate configurations (Hsu & Kakade, 2012; Goyal et al., 2014; Bhaskara et al., 2014a; Anderson et al., 2014; Ge et al., 2015).

Algorithmic approaches based on method-of-moments and tensor decompositions rely heavily on the exact parametric form of the Gaussian distribution and the exact algebraic expressions to express various moments of the distribution in terms of the parameters. These algebraic methods can be easily foiled by a monotone adversary, since the adversary can perturb any subset to alter the moments significantly (for example, even the first moment, i.e., the mean of a cluster, can change by  $\Omega(\sigma)$ ).

Recent work has also focused on provable guarantees for heuristics such as Maximum Likelihood estimation and the Expectation Maximization (EM) algorithm for parameter estimation (Dasgupta & Schulman, 2007; Balakrishnan et al., 2014; Xu et al., 2016; Daskalakis et al., 2016; Tang & Monteleoni, 2017). Very recently, (Regev & Vijayaraghavan, 2017) considered other iterative algorithms for parameter estimation of spherical Gaussians, and studied the optimal order of separation required for parameter estimation. However, we are not aware of any existing analysis that shows that these iterative algorithms for parameter estimation are robust to modeling errors.

Another recent line of exciting work concerns designing robust high-dimensional estimators of the mean and covariance of a single Gaussian (and mixtures of  $k$  Gaussians) when an  $\varepsilon = \Omega_k(1)$  fraction of the points are adversarially corrupted (Diakonikolas et al., 2016; Lai et al., 2016; Charikar et al., 2017). However, these results and similar results on agnostic learning do not necessarily recover the ground-truth clustering. Further, they typically assume that only a  $o(1/k)$  fraction of the points are corrupted, while potentially all the points could be perturbed in the semi-random model. On the other hand, our work does not necessarily give guarantees for estimating the means of the original Gaussians (in fact the centers given by the planted clustering in the semi-random instance can be  $\Omega(\sigma)$  far from the original means). Hence, our semi-random model is incomparable to the model of robustness considered in these works.

Finally in concurrent and independent works (Hopkins & Li, 2017; Kothari & Steinhardt, 2017; Diakoniko-

las et al., 2017), algorithms based on sum-of-squares relaxations and robust estimation techniques were used to obtain algorithms with run time (and sample-complexity) of  $(dk)^{O(1/\varepsilon)}$  for clustering mixtures of spherical Gaussians when the separation is of the order of  $k^\varepsilon$  for any constant  $\varepsilon > 0$ . Some of these results also tolerate a small fraction of each cluster containing outliers as in (Diakonikolas et al., 2016). However, to the best of our knowledge these guarantees does not work in our semi-random model, and the results are incomparable like the above works on robust estimation (Diakonikolas et al., 2016; Lai et al., 2016; Charikar et al., 2017). Further, our algorithmic guarantees are for Lloyd’s algorithm and  $k$ -means++, which are practical heuristics that form the method-of-choice in practice.

**Semi-random models for other optimization problems.** There has been a long line of work on the study of semi-random models for various optimization problems. Blum and Spencer (Blum & Spencer, 1995) initiated the study of semi-random models, and studied the problem of graph coloring. Feige and Kilian (Feige & Kilian, 1998) considered semi-random models involving monotone adversaries for various problems including graph partitioning, independent set and clique. Makarychev et al. (Makarychev et al., 2012; 2014) designed algorithms for more general semi-random models for various graph partitioning problems. The work of (Moitra et al., 2015) studied the power of monotone adversaries in the context of community detection (stochastic block models), while (Makarychev et al., 2016) considered the robustness of community detection to monotone adversaries and different kinds of errors and model misspecification. Semi-random models have also been studied for correlation clustering (Mathieu & Schudy, 2010; Makarychev et al., 2015), noisy sorting (Makarychev et al., 2013) and coloring (David & Feige, 2016).

## 7. Preliminaries and Semi-random model

We first formally define the Gaussian mixture model.

**Definition 7.1.** (Gaussian Mixture Model). A Gaussian mixture model with  $k$  components is defined by the parameters  $(\mu_1, \mu_2, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k, w_1, \dots, w_k)$ . Here  $\mu_i \in \mathbb{R}^d$  is the mean for component  $i$  and  $\Sigma_i \in \mathbb{S}_+^d$  is the corresponding  $d \times d$  covariance matrix.  $w_i \in [0, 1]$  is the mixing weight and we have that  $\sum_{i=1}^k w_i = 1$ . An instance  $\mathcal{X} = \{x^{(1)}, \dots, x^{(N)}\}$  from the mixture is generated as follows: for each  $t \in [N]$ , sample a component  $i \in [k]$  independently at random with prob-

ability  $w_i$ . Given the component, sample  $x^{(t)}$  from  $\mathcal{N}(\mu_i, \Sigma_i)$ . The  $N$  points can be naturally partitioned into  $k$  clusters  $C_1, \dots, C_k$  where cluster  $C_i$  corresponds to the points that are sampled from component  $i$ . We will refer to this as the *planted clustering* or *ground truth clustering*.

Clustering data from a mixture of Gaussians is a natural average-case model for the  $k$ -means clustering problem. Specifically, if the means of a Gaussian mixture model are well separated, then with high probability, the ground truth clustering of an instance sampled from the model corresponds to the  $k$ -means optimal clustering.

**Definition 7.2.** ( $k$ -means clustering). Given an instance  $\mathcal{X} = \{x^{(1)}, \dots, x^{(N)}\}$  of  $N$  points in  $\mathbb{R}^d$ , the  $k$ -means problem is to find  $k$  points  $\mu_1, \dots, \mu_k$  such as to minimize  $\sum_{t \in [N]} \min_{i \in [k]} \|x^{(t)} - \mu_i\|^2$ .

The optimal means or centers  $\mu_1, \dots, \mu_k$  naturally define a clustering of the data where each point is assigned to its closest cluster. A key property of the  $k$ -means objective is that the optimal solution induces a locally optimal clustering.

**Definition 7.3.** (Locally Optimal Clustering). A clustering  $C_1, \dots, C_k$  of  $N$  data points in  $\mathbb{R}^d$  is locally optimal if for each  $i \in [k]$ ,  $x^{(t)} \in C_i$ , and  $j \neq i$  we have that  $\|x^{(t)} - \mu(C_i)\| \leq \|x^{(t)} - \mu_j\|$ . Here  $\mu(C_i)$  is the average of the points in  $C_i$ .

Hence, given the optimal  $k$ -means clustering, the optimal centers can be recovered by simply computing the average of each cluster. This is the underlying principle behind the popular Lloyd’s algorithm (Lloyd, 1982) for  $k$ -means clustering. The algorithm starts with a choice of initial centers. It then repeatedly computes new centers to be the average of the clusters induced by the current centers. Hence the algorithm converges to a locally optimal clustering. Although popular in practice, the worst case performance of Lloyd’s algorithm can be arbitrarily bad (Arthur & Vassilvitskii, 2005). The choice of initial centers is very important in the success of the Lloyd’s algorithm. We show that our theoretical guarantees hold when the initialization is done via the popular  $k$ -means++ algorithm (Arthur & Vassilvitskii, 2007). There also exist more sophisticated constant factor approximation algorithms for the  $k$ -means problem (Kanungo et al., 2002; Ahmadian et al., 2016) that can be used for seeding in our framework.

While the clustering  $C_1, C_2, \dots, C_k$  typically represents a partition of the index set  $[N]$ , we will sometimes abuse notation and use  $C_i$  to also denote the set of points in  $\mathcal{X}$  that correspond to these indices in  $C_i$ . Finally, many of the statements are probabilistic in nature depending on the randomness in the semi-random model. In the

following section, w.h.p. will refer to a probability of at least  $1 - o(1)$  (say  $1 - 1/\text{poly}(N)$ ), unless specified otherwise.

### 7.1. Properties of Semi-random Gaussians

In this section we state and prove properties of semi-random mixtures that will be used throughout the analysis in the subsequent sections. We first start with a couple of simple lemmas that follow directly from the corresponding lemmas about high dimensional Gaussians.

**Lemma 7.4.** *Consider any semi-random instance  $\mathcal{X} = \{x^{(1)}, \dots, x^{(N)}\}$  with parameters  $\mu_1, \dots, \mu_k, \sigma^2$  and clusters  $C_1, \dots, C_k$ . Then with high probability we have*

$$\forall i \in [k], \forall \ell \in C_i, \quad \|x^{(\ell)} - \mu_i\|_2 \leq \sigma(\sqrt{d} + 2\sqrt{\log N}). \quad (9)$$

*Proof.* Let  $y^{(t)}$  denote the point generated in the semi-random model in step 2 (Definition 1.1) before the semi-random perturbation was applied. Let  $\bar{x}^{(t)} = x - \mu_i$ ,  $\bar{y}^{(t)} = y - \mu_i$  where  $t \in C_i$ . We have

$$\forall i \in [k], \forall t \in C_i, \quad \|\bar{x}^{(t)}\|_2 \leq \|\bar{y}^{(t)}\|_2 \leq \sigma(\sqrt{d} + 2\sqrt{\log N}),$$

from Lemma A.3.  $\square$

**Lemma 7.5.** *Consider any semi-random instance  $\mathcal{X} = \{x^{(1)}, \dots, x^{(N)}\}$  with parameters  $\mu_1, \dots, \mu_k, \sigma^2$  and clusters  $C_1, \dots, C_k$ , and let  $u$  be a fixed unit vector in  $\mathbb{R}^d$ . Then with probability at least  $(1 - 1/(N^3))$  we have*

$$\forall i \in [k], t \in C_i, \quad |\langle x^{(t)} - \mu_i, u \rangle| < 3\sigma\sqrt{\log N}. \quad (10)$$

*Proof.* Let  $y^{(t)}$  denote the point generated in the semi-random model in step 2 (Definition 1.1) before the semi-random perturbation was applied. Let  $\bar{x}^{(t)} = x - \mu_i$ ,  $\bar{y}^{(t)} = y - \mu_i$  where  $t \in C_i$ .

Consider the sample  $t \in C_i$ . Let  $\Sigma_i$  be the covariance matrix of  $i$ th Gaussian component; hence  $\|\Sigma_i\| \leq \sigma$ . The projection  $\langle \bar{y}^{(t)}, u \rangle$  is a Gaussian with mean 0 and variance  $u^T \Sigma_i u \leq \sigma^2$ . From Lemma A.1

$$\begin{aligned} \mathbb{P} \left[ |\langle \bar{x}^{(t)}, u \rangle| \geq 3\sigma\sqrt{\log N} \right] &\leq \mathbb{P} \left[ |\langle \bar{y}^{(t)}, u \rangle| \geq 3\sigma\sqrt{\log N} \right] \\ &\leq \exp(-4 \log N) \leq N^{-4}. \end{aligned}$$

Hence from a union bound over all  $N$  samples, the lemma follows.  $\square$

The above lemma immediately implies the following lemma after a union bound over the  $k^2 < N^2$  directions given by the unit vectors along  $(\mu_i - \mu_j)$  directions.

**Lemma 7.6.** *Consider any semi-random instance  $\mathcal{X} = \{x^{(1)}, \dots, x^{(N)}\}$  with parameters  $\mu_1, \dots, \mu_k, \sigma^2$  and clusters  $C_1, \dots, C_k$ . Then with high probability we have*

$$\forall i \in [k], t \in C_i, \quad \left| \left\langle x^{(t)} - \mu_i, \frac{\mu_i - \mu_j}{\|\mu_i - \mu_j\|_2} \right\rangle \right| < 3\sigma\sqrt{\log N}. \quad (11)$$

We next state a lemma about how far the mean of the points in a component of a semi-random GMM can move away from the true parameters.

**Lemma 7.7.** *Consider any semi-random instance  $\mathcal{X}$  with  $N$  points generated with parameters  $\mu_1, \dots, \mu_k, C_1, \dots, C_k$  such that  $N_i \geq 4(d + \log(\frac{k}{\delta}))$  for all  $i \in [k]$ . Then with probability at least  $1 - \delta$  we have that*

$$\forall i \in [k], \left\| \frac{1}{|C_i|} \sum_{x \in C_i} x - \mu_i \right\|_2 \leq 2\sigma. \quad (12)$$

*Proof.* For each point  $x \in C_i$  in the semi-random GMM, let  $y_x$  be the original point in the GMM that is modified to produce  $x$ . Then, we know that  $x - \mu_i = \lambda_x(y_x - \mu_i)$  where  $\lambda_x \in [0, 1]$ . Hence,  $\frac{1}{|C_i|} \sum_{x \in C_i} (x - \mu_i) = \frac{1}{|C_i|} A_i D v$ , where  $A_i$  is the matrix with columns as  $(y_x - \mu_i)$  for  $x \in C_i$ ,  $D$  is a diagonal matrix with values  $\lambda_x$ , and  $v$  is a unit length vector in the direction of  $\frac{1}{|C_i|} \sum_{x \in C_i} (x - \mu_i)$ . Then, we have that  $\left\| \frac{1}{|C_i|} \sum_{x \in C_i} x - \mu_i \right\| = \left\| \frac{1}{|C_i|} A_i D v \right\| \leq \frac{1}{|C_i|} \|A\| \leq 2\sigma$  (from A.5).  $\square$

The next lemma argues about the variance of component  $i$  around  $\mu_i$  in a semi-random GMM.

**Lemma 7.8.** *Consider any semi-random instance  $\mathcal{X}$  with  $N$  points generated with parameters  $\mu_1, \dots, \mu_k, C_1, \dots, C_k$  such that  $N_i \geq 4(d + \log(\frac{k}{\delta}))$  for all  $i \in [k]$ . Then with probability at least  $1 - \delta$  we have that*

$$\forall i \in [k], \max_{v: \|v\|=1} \frac{1}{|C_i|} \sum_{x \in C_i} |\langle x - \mu_i, v \rangle|^2 \leq 4\sigma^2. \quad (13)$$

*Proof.* Exactly as in the proof of Lemma 7.7, we can write  $\max_{v: \|v\|=1} \frac{1}{|C_i|} \sum_{x \in C_i} |\langle x - \mu_i, v \rangle|^2 = \max_{v: \|v\|=1} \frac{1}{|C_i|} \sum_{x \in C_i} |\lambda_x^2 \langle y_x - \mu_i, v \rangle|^2 \leq \max_{v: \|v\|=1} \frac{1}{|C_i|} \sum_{x \in C_i} |\langle y_x - \mu_i, v \rangle|^2$ . Furthermore, since  $y_x$  are points from a Gaussian we know that with probability at least  $1 - \delta$ , for all  $i \in [k]$ ,  $\max_{v: \|v\|=1} \frac{1}{|C_i|} \sum_{x \in C_i} |\langle y_x - \mu_i, v \rangle|^2 \leq 4\sigma^2$ . Hence, the claim follows.  $\square$

We would also need to argue about the mean of a large subset of points from a component of a semi-random GMM.

**Lemma 7.9.** *Consider any semi-random instance  $\mathcal{X}$  with  $N$  points generated with parameters  $\mu_1, \dots, \mu_k$  and planted clustering  $C_1, \dots, C_k$  such that  $N_i \geq 16(d + \log(\frac{k}{\delta}))$  for all  $i \in [k]$ . Let  $G_i \subseteq C_i$  be such that  $|G_i| \geq (1 - \varepsilon)|C_i|$  where  $\varepsilon < \frac{1}{2}$ . Then, with probability at least  $1 - \delta$ , we have that*

$$\forall i \in [k], \|\mu(G_i) - \mu_i\| \leq (4 + \frac{2}{\sqrt{1 - \varepsilon}})\sigma. \quad (14)$$

*Proof.* Let  $C_i$  be the set of points in component  $i$  and let  $\nu_i$  be the mean of the points in  $C_i$ . Notice that from Lemma 7.7 and the fact that the perturbation is semi-random, we have that with probability at least  $1 - \frac{\delta}{2}$ ,  $\|\nu_i - \mu_i\| \leq 2\sigma$ . Also, because the component is a semi-random perturbation of a Gaussian, we have from Lemma 7.8 that  $\frac{1}{|C_i|} \max_{v: \|v\|=1} \sum_{x \in C_i} \langle x - \nu_i, v \rangle^2 \leq 4\sigma^2$  with probability at least  $1 - \frac{\delta}{2}$ .

Hence, with probability at least  $1 - \delta$  we have that  $\|\mu(G_i) - \mu_i\| \leq \|\nu_i - \mu_i\| + \|\mu(G_i) - \nu_i\| \leq 4\sigma + \|\mu(G_i) - \nu_i\|$ . To bound the second term notice that  $\|\mu(G_i) - \nu_i\| = |(\frac{1}{|G_i|} \sum_{x \in G_i} \langle x - \nu_i, \hat{u} \rangle)|$ , where  $\hat{u}$  is a unit vector in the direction of  $(\mu(G_i) - \nu_i)$ . Using Cauchy-Schwarz inequality, this is at most  $\frac{1}{\sqrt{|G_i|}} \sqrt{\sum_{x \in G_i} \langle x - \nu_i, \hat{u} \rangle^2} \leq \frac{2\sigma}{\sqrt{1 - \varepsilon}}$ . Combining the two bounds gives us the result.  $\square$

Finally, we argue about the variance of the entire data matrix of a semi-random GMM.

**Lemma 7.10.** *Consider any semi-random instance  $\mathcal{X}$  with  $N$  points generated with parameters  $\mu_1, \dots, \mu_k, C_1, \dots, C_k$  such that  $N_i \geq 4(d + \log(\frac{k}{\delta}))$  for all  $i \in [k]$ . Let  $A \in \mathbb{R}^{d \times N}$  be the matrix of data points and let  $M \in \mathbb{R}^{d \times N}$  be the matrix composed of the means of the corresponding clusters. Then, with probability at least  $1 - \delta$ , we have that*

$$\|A - M\| \leq 4\sigma\sqrt{N}. \quad (15)$$

*Proof.* Let  $M^*$  be the matrix of true means corresponding to the cluster memberships. We can write  $\|A - M\| \leq \|A - M^*\| + \|M^* - M\|$ . Using Lemma 7.7, we know that with probability at least  $1 - \frac{\delta}{2}$ ,  $\max_i \|M_i^* - M_i\| \leq 2\sigma$ . Hence,  $\|M^* - M\| \leq 2\sigma\sqrt{N}$ . Furthermore,  $\|A - M^*\|^2 = \max_{v: \|v\|=1} \sum_i \sum_{x \in C_i} |(x - \mu_i) \cdot v|^2$ . From Lemma 7.8, with probability at least  $1 - \frac{\delta}{2}$ , we can bound the sum by at most  $4\sigma^2 N$ . Hence,  $\|A - M^*\| \leq 2\sigma\sqrt{N}$ . Combining the two bounds we get the claim.  $\square$

## 8. Upper Bounds for Semi-random GMMs

In this section we prove the following theorem that provides algorithmic guarantees for the Lloyd's algorithm with appropriate initialization, under the semi-random model for mixtures of Gaussians in Definition 1.1.

**Theorem 8.1.** *There exists a universal constant  $c_0, c_1 > 0$  such that the following holds. There exists a polynomial time algorithm that for any semi-random instance  $\mathcal{X}$  on  $N$  points with planted clustering  $C_1, \dots, C_k$  generated by the semi-random GMM model (Definition 1.1) with parameters  $\mu_1, \dots, \mu_k, \sigma^2$  s.t.*

$$\forall i \neq j \in [k], \|\mu_i - \mu_j\|_2 > \Delta\sigma \quad (16)$$

where  $\Delta > c_0 \sqrt{\min\{k, d\} \log N}$  and  $N \geq k^2 d^2 / w_{min}^2$  finds w.h.p. a clustering  $C'_1, C'_2, \dots, C'_k$  such that

$$\min_{\pi \in \text{Perm}_k} \sum_{i=1}^k |C_{\pi(i)} \Delta C'_i| \leq \frac{c_1 k d}{\Delta^4} \cdot \max \left\{ 1, \log \left( \frac{3(\sqrt{d} + 2)\sqrt{\log N}}{\Delta^2} \right) \right\}.$$

In Section 9 we show that the above error bound is close to the information theoretically optimal bound (up to the logarithmic factor). The Lloyd's algorithm as described in Figure 1 consists of two stages, the initialization stage and an iterative improvement stage.

1. Let  $A$  be the  $N \times d$  data matrix with rows  $A_i$  for  $i \in [N]$ . Use  $A$  to compute initial centers  $\mu_0^{(1)}, \mu_0^{(2)}, \dots, \mu_0^{(k)}$  as detailed in Proposition 8.2.
2. Use these  $k$ -centers to seed a series of Lloyd-type iterations. That is, for  $r = 1, 2, \dots$  do:
  - Set  $Z_i$  be the set of points for which the closest center among  $\mu_{r-1}^{(1)}, \mu_{r-1}^{(2)}, \dots, \mu_{r-1}^{(k)}$  is  $\mu_{r-1}^{(i)}$ .
  - Set  $\mu_r^{(i)} \leftarrow \frac{1}{|Z_i|} \sum_{A_j \in Z_i} A_j$ .

Figure 1. Lloyd's Algorithm

The initialization follows the same scheme as proposed by Kumar and Kannan in (Kumar & Kannan, 2010). The initialization algorithm first performs a  $k$ -SVD of the data matrix followed by running the  $k$ -means++ algorithm (Arthur & Vassilvitskii, 2007) that uses  $D^2$ -sampling to compute seed centers. One can also use any constant factor approximation algorithm for  $k$ -means clustering in the projected space to obtain the initial centers (Kanungo et al., 2002; Ahmadian et al.,

2016). This approach works for clusters that are nearly balanced in size. However, when the cluster sizes are arbitrary, an appropriate transformation of the data is performed first that amplifies the separation between the centers. Following this transformation, the ( $k$ -SVD +  $k$ -means++) is used to get the initial centers. The formal guarantee of the initialization procedure is encapsulated in the following proposition, whose proof is given in Section 8.2.

The main algorithmic contribution of this paper is an analysis of the Lloyd's algorithm when the points come from the semi-random GMM model. For the rest of the analysis we will assume that the instance  $\mathcal{X}$  generated from the semi-random GMM model satisfies (9) to (17). These eight equations are shown to hold w.h.p. in Section 7.1 for instances generated from the model. Our analysis will in fact hold for any deterministic data set satisfying these equations. This helps to gracefully argue about performing many iterations of Lloyd's on the same data set without the need to draw fresh samples at each step.

**Proposition 8.2.** *In the above notation for any  $\delta > 0$ , suppose we are given an instance  $\mathcal{X}$  on  $N$  points satisfying (9)-(17) such that  $|C_i| \geq \Omega(d + \log(\frac{k}{\delta}))$  and assume that  $\Delta \geq 125\sqrt{\min\{k, d\} \log N}$ . Then after the initialization step, for every  $\mu_i$  there exists  $\mu'_i$  such that  $\|\mu_i - \mu'_i\| \leq \tau\sigma$ , where  $\tau < \Delta/24$ .*

The analysis of the Lloyd's iterations crucially relies on the following lemma that upper bounds the number of misclassified points when the current Lloyd's iterative is relatively close to the true means.

**Lemma 8.3** (Projection condition). *In the above notation, consider an instance  $\mathcal{X}$  satisfying (9)-(17) and (16) and suppose we are given  $\mu'_1, \dots, \mu'_k$  satisfying  $\forall j \in [k], \|\mu'_j - \mu_j\|_2 \leq \tau\sigma$  and  $\tau < \Delta/24$ . Then there exists a set  $Z \subset \mathcal{X}$  such that for any  $i \in [k]$  we have*

$$\forall x \in C_i \cap (\mathcal{X} \setminus Z), \quad \|x - \mu'_i\|_2^2 \leq \min_{j \neq i} \|x - \mu'_j\|_2^2, \quad \text{where}$$

$$|Z| = O\left(\frac{d\tau^2}{\Delta^4} \cdot \max\left\{1, \log\left(\frac{3\tau(\sqrt{d}+2\sqrt{\log N})}{\Delta^2}\right)\right\}\right).$$

The following lemma quantifies the improvement in each step of the Lloyd's algorithm. The proof uses Lemma 8.3 along with properties of semi-random Gaussians.

**Lemma 8.4.** *In the above notation, suppose we are given an instance  $\mathcal{X}$  on  $N$  points with  $w_i N \geq \frac{d\sqrt{d}}{4\log(d)}$  for all  $i$  satisfying (9)-(17). Furthermore, suppose we are given centers  $\mu'_1, \dots, \mu'_k$  such that  $\|\mu'_i - \mu_i\| \leq \tau\sigma$ ,  $\forall i \in [k]$  where  $\tau < \Delta/24$ . Then the centers  $\mu''_1, \dots, \mu''_k$  obtained after one Lloyd's update satisfy  $\|\mu''_i - \mu_i\| \leq \max((6 + \frac{\tau}{4})\sigma, \frac{\tau}{2}\sigma)$  for all  $i \in [k]$ .*

We now present the proof of Theorem 8.1.

*Proof of Theorem 8.1.* Firstly, the eight deterministic conditions (9)-(17) are shown to hold for instance  $\mathcal{X}$  w.h.p. in Section 7.1. The proof follows in a straightforward manner by combining Proposition 8.2, Lemma 8.4 and Lemma 8.3. Proposition 8.2 shows that  $\|\mu_i^{(0)} - \mu_i\|_2 \leq \Delta/(24)$  for all  $i \in [k]$ . Applying Lemma 8.4, we have that after  $T = O(\log \Delta)$  iterations we get  $\|\mu_i^{(T)} - \mu_i\|_2 \leq 8\sigma$  for all  $i \in [k]$  w.h.p. Finally using Lemma 8.3 with  $\tau = 1$ , the theorem follows.  $\square$

### 8.1. Analyzing Lloyd's Algorithm

The following lemma is crucial in analyzing the performance of the Lloyd's algorithm. We would like to upper bound the inner product  $|\langle x^{(\ell)} - \mu_i, \hat{e} \rangle| < \lambda\sigma$  for every direction  $\hat{e}$  and sample  $\ell \in [N]$ , but this is impossible since  $\hat{e}$  can be aligned along  $x^{(\ell)} - \mu_i$ . The following lemma however upper bounds the total number of points in the dataset that can have a large projection of  $\lambda$  (or above) onto any direction  $\hat{e}$  by at most  $\tilde{O}(d/\lambda^2)$ . This involves a union bound over a net of all possible directions  $\hat{e}$ .

**Lemma 8.5** (Points in Bad Directions). *Consider any semi-random instance  $\mathcal{X} = \{x^{(1)}, \dots, x^{(N)}\}$  with  $N$  points having parameters  $\mu_1, \dots, \mu_k, \sigma^2$  and planted clustering  $C_1, \dots, C_k$ , and suppose  $\forall i \in [k], \ell \in C_i, \bar{x}^{(\ell)} = x^{(\ell)} - \mu_i$ . Then there exists a universal constant  $c > 0$  s.t. for any  $\lambda > 100\sqrt{\log N}$ , with probability at least  $1 - 2^{-d}$ , we have that  $\forall \hat{e} \in \mathbb{R}^d$  s.t.  $\|\hat{e}\|_2 = 1$ ,*

$$\left| \left\{ \ell \in [N] : |\langle \bar{x}^{(\ell)}, \hat{e} \rangle| > \lambda\sigma \right\} \right| \leq \frac{cd}{\lambda^2} \cdot \max\left\{1, \log\left(\frac{3(\sqrt{d}+2\sqrt{\log N})}{\lambda}\right)\right\}. \quad (17)$$

*Proof.* Set  $\eta := \min\{\lambda/(2\sqrt{d} + 2\sqrt{\log N}), \frac{1}{2}\}$  and  $m := 512d \log(3/\eta)/\lambda^2$ . Consider an  $\eta$ -net  $\mathcal{N} \subset \{u : \|u\|_2 = 1\}$  over unit vectors in  $\mathbb{R}^d$ . Hence

$$\forall u \in \mathbb{R}^d : \|u\|_2 = 1, \exists v \in \mathcal{N} \text{ s.t. } \|u - v\|_2 \leq \eta \text{ and } |\mathcal{N}| \leq \left(\frac{2+\eta}{\eta}\right)^d \leq \exp(d \log(3/\eta)).$$

Further, since  $|\langle \bar{x}, \hat{e} \rangle| > \lambda$  and  $\mathcal{N}$  is an  $\eta$ -net, there exists some unit vector  $u = u(\hat{e}) \in \mathcal{N}$

$$|\langle \bar{x}, u \rangle| > |\langle \bar{x}, \hat{e} \rangle + \langle \bar{x}, \hat{e} - u \rangle| \geq \sigma\lambda - \|\bar{x}\|_2 \|\hat{e} - u\|_2 \quad (18)$$

$$\geq \sigma(\lambda - \eta(\sqrt{d} + 2\sqrt{\log N})) \geq \frac{\lambda}{2}, \quad (19)$$

Consider a fixed  $x \in \{x^{(1)}, \dots, x^{(N)}\}$  and a fixed direction  $u \in \mathcal{N}$ . Since the variance of  $y$  is at most  $\sigma^2$

we have

$$\mathbb{P} \left[ |\langle \bar{x}, u \rangle| > \lambda\sigma/2 \right] \leq \mathbb{P} \left[ |\langle \bar{y}, u \rangle| > \lambda\sigma/2 \right] \leq \exp(-\lambda^2/8).$$

The probability that  $m$  points in  $\{x^{(1)}, \dots, x^{(N)}\}$  satisfy (19) for a fixed direction  $u$  is at most  $\binom{N}{m} \cdot \exp(-m\gamma^2/2)$ . Let  $E$  represent the bad event that there exists a direction in  $\mathcal{N}$  such that more than  $m$  points satisfy the bad event given by (19). The probability of  $E$  is at most

$$\begin{aligned} \mathbb{P}[E] &\leq |\mathcal{N}| \cdot \binom{N}{m} \exp(-m\lambda^2/8) \\ &\leq \exp\left(d \log(3/\eta) + m \log N - \frac{m\lambda^2}{8}\right) \\ &\leq \exp\left(-d \log(1/\eta)\right) \leq \eta^d, \end{aligned}$$

since for our choice of parameters  $\lambda^2 > 32 \log N$ , and  $m\lambda^2 \geq 32d \log(3/\eta)$ .

□

We now analyze each iteration of the Lloyd's algorithm and show that we make progress in each step by misclassifying fewer points with successive iterations. As a first step we begin with the proof of Lemma 8.3.

*Proof of Lemma 8.3.* Set  $m := 512d \log(3/\eta) \tau^2 / \Delta^4$  where  $\eta = \min \left\{ \Delta^2 / (\tau(2\sqrt{d} + 2\sqrt{\log N})), \frac{1}{2} \right\}$ .

Fix a sample  $x \in \{x^{(1)}, \dots, x^{(N)}\}$  and suppose  $x \in C_i$  and let  $y := y(x)$  be the corresponding point before the semi-random perturbation, and let  $\bar{x} = x - \mu_i$ ,  $\bar{y} = y - \mu_i$ . For each  $i \in [k]$ , let  $\hat{e}_i$  be the unit vector along  $(\mu_i - \mu'_i)$ .

We first observe that by projecting the Gaussians around  $\mu_i, \mu_j$  onto the direction along  $\hat{e}_{ij} = (\mu_i - \mu_j) / \|\mu_i - \mu_j\|_2$ , we have that

$$\begin{aligned} \|x - \mu_j\|_2^2 - \|x - \mu_i\|_2^2 &\geq \langle x - \mu_j, \hat{e}_{ij} \rangle^2 - \langle x - \mu_i, \hat{e}_{ij} \rangle^2 \\ &\geq (|\langle x - \mu_j, \hat{e}_{ij} \rangle| - |\langle x - \mu_i, \hat{e}_{ij} \rangle|)^2 \\ &\geq (|\langle \mu_i - \mu_j, \hat{e}_{ij} \rangle| - 2|\langle x - \mu_i, \hat{e}_{ij} \rangle|)^2 \\ &\geq (\Delta\sigma - 2|\langle x - \mu_i, \hat{e}_{ij} \rangle|)^2 \\ &\geq (\Delta\sigma - 6\sigma\sqrt{\log N})^2 \geq \frac{1}{4}\Delta^2\sigma^2, \end{aligned} \tag{20}$$

where the first inequality follows from (11), and the second inequality uses  $\Delta > 12\sqrt{\log N}$ .

Suppose  $x \in C_i$  is misclassified i.e.,  $\|x - \mu'_j\|_2 \geq \|x - \mu_j\|_2$  for some  $j \in [k] \setminus \{i\}$ . Then applying triangle

inequality and rearranging we get,

$$\begin{aligned} &\left| \left\langle \bar{x}, \frac{\mu_i - \mu'_i}{\|\mu_i - \mu'_i\|_2} \right\rangle \right| + \left| \left\langle \bar{x}, \frac{\mu_j - \mu'_j}{\|\mu_j - \mu'_j\|_2} \right\rangle \right| \\ &\geq \frac{\left(\frac{\Delta^2}{8} - \frac{\tau^2}{2} - \tau\Delta\right)\sigma^2}{\tau\sigma} \\ &\geq \frac{\Delta^2}{16\tau}\sigma, \end{aligned}$$

since  $\tau < \Delta/(24)$ . Hence, we have that if  $x \in C_i$  is misclassified by  $\mu'_1, \dots, \mu'_k$  then

$$|\langle \bar{x}, \hat{e} \rangle| > \sigma\Delta^2/(32\tau) \text{ for some unit vector } \hat{e} \in \mathbb{R}^d. \tag{21}$$

From (17) with  $\lambda = \Delta^2/(32\tau)$ , we get from (17) that at most  $m$  points in  $C_i$  can satisfy (21). Hence the lemma follows. □

Next we prove Lemma 8.4, which quantifies the improvement in every iteration of the Lloyd's algorithm.

*Proof of Lemma 8.4.* Let  $C_1, C_2, \dots, C_k$  be the partitioning of the indices according to the ground truth clustering of the semi-random instance  $\mathcal{X}$  and  $S_1, S_2, \dots, S_k$  be the indices of the clustering obtained by using the centers  $\mu'_i$ . Then  $\mu''_i = \frac{1}{|S_i|} \sum_{t \in S_i} x^{(t)}$ . Partition  $S_i$  into two sets  $G_i$  and  $B_i$  where  $G_i = S_i \cap C_i$  and  $B_i = S_i \setminus G_i$ . Let  $\mu(G_i)$  and  $\mu(B_i)$  be the means of the two partitions respectively.

Let  $\gamma = O\left(\frac{d\tau^2}{\Delta^4} \max\left\{1, \log\left(\frac{3\tau(\sqrt{d}+2\sqrt{\log N})}{\Delta^2}\right)\right\}\right)$ . From Lemma 8.3 we know that  $|G_i| \geq |C_i| - \gamma$  and  $|B_i| \leq k\gamma$ . Then we have that  $\mu''_i = \frac{|G_i|}{|S_i|} \mu(G_i) + \frac{|B_i|}{|S_i|} \mu(B_i)$ . Hence,

$\|\mu''_i - \mu_i\| \leq \frac{|G_i|}{|S_i|} \|\mu(G_i) - \mu_i\| + \frac{|B_i|}{|S_i|} \|\mu(B_i) - \mu_i\|$ . We have  $\frac{|G_i|}{|C_i|} \geq 1 - \frac{\gamma}{|C_i|} \geq 1 - \frac{\tau}{64\sqrt{k}\sqrt{d}}$  using the bound on  $\Delta$  and  $|C_i| = w_i N \geq \frac{d\sqrt{d}}{4\log(d)}$ . Using (14) we get that

$$\begin{aligned} \frac{|G_i|}{|S_i|} \|\mu(G_i) - \mu_i\| &\leq \left(4 + \frac{2}{\sqrt{1 - \frac{\tau}{64\sqrt{k}\sqrt{d}}}}\right)\sigma \\ &\leq \left(6 + \frac{\tau}{128\sqrt{k}\sqrt{d}}\right)\sigma \\ &\leq 6\sigma + \frac{\tau}{8}\sigma. \end{aligned}$$

To bound the second term we first show that for each point  $x^{(t)} \in B_i$ ,  $\|x^{(t)} - \mu_i\| \leq (\sqrt{d} + 2\sqrt{\log N} + 2\tau)\sigma$ . Let  $C_j$  be the cluster that point  $x^{(t)}$  belongs to. Then

$$\begin{aligned} \|x^{(t)} - \mu_i\| &\leq \|x^{(t)} - \mu'_j\| + \tau\sigma \leq \|x^{(t)} - \mu'_j\| + \tau\sigma \\ &\leq \|x^{(t)} - \mu_j\| + 2\tau\sigma \\ &\leq (\sqrt{d} + 2\sqrt{\log N} + 2\tau)\sigma, \end{aligned}$$

using (9). Hence,

$$\begin{aligned} \frac{|B_i|}{|S_i|} \|\mu(B_i) - \mu_i\| &\leq \frac{|B_i|}{|S_i|} (\sigma\sqrt{d} + \sigma\sqrt{\log N} + 2\tau\sigma) \\ &\leq \frac{2k\gamma}{|C_i|} (\sigma\sqrt{d} + \sigma\sqrt{\log N} + 2\tau\sigma) \\ &< \frac{\tau}{8}\sigma. \end{aligned}$$

Combining, we get that  $\|\mu'_i - \mu_i\| \leq (6 + \frac{\tau}{4})\sigma \leq \max(6\sigma + \frac{\tau}{4}, \frac{\tau}{2}\sigma)$ .  $\square$

## 8.2. Initialization

In this section we describe how to obtain the initial centers satisfying the condition in Lemma 8.4. The final initialization procedure relies on the following subroutine that provides a good initializer if the mean separation is much larger than that in Theorem 8.1. Let  $A$  denote the  $N \times d$  matrix of data points and  $M^*$  be the  $N \times d$  matrix where each row of  $C$  is equal to one of the means  $\mu_i$ s of the component to which the corresponding row of  $A$  belongs to.

**Lemma 8.6.** *In the above notation, for any  $\delta > 0$  suppose we are given an instance  $\mathcal{X}$  on  $N$  points satisfying (9)-(17), with components  $C_1, \dots, C_k$  such that  $|C_i| \geq \Omega(d + \log(\frac{k}{\delta}))$ . Let  $A$  be the  $N \times d$  matrix of data points and  $\hat{A}$  be the matrix obtained by projecting points onto the best  $k$ -dimensional subspace obtained by SVD of  $A$ . Let  $\mu'_i$  be the centers obtained by running an  $\alpha$  factor  $k$ -means approximation algorithm on  $A$ . Then for every  $\mu_i$  there exists  $\mu'_i$  such that  $\|\mu_i - \mu'_i\| \leq 20\sqrt{k}\alpha \frac{\|A - M^*\|}{\sqrt{Nw_{\min}}}$ .*

*Proof.* Let  $\hat{A}$  denote the matrix obtained by projecting  $A$  onto the span of its top  $k$  right singular vectors. Furthermore, let  $\nu_1, \dots, \nu_k$  be the centers obtained by running a 9-approximation algorithm for  $k$ -means on the instance  $\hat{A}$ . We know that the optimal  $k$ -means solution for  $\hat{A}$  is at most  $\|\hat{A} - M^*\|_F^2$ . Since both  $\hat{A}$  and  $M^*$  are rank  $k$  matrices, we get that  $\|\hat{A} - M^*\|_F^2 \leq 2k\|\hat{A} - M^*\|_2^2 \leq 2k(\|\hat{A} - A\|_2^2 + \|A - M^*\|_2^2)$ . Since  $\hat{A}$  is the best rank  $k$  approximation to  $A$  we also have that  $\|\hat{A} - A\|_2^2 \leq \|A - M^*\|_2^2$ . Hence,  $\|\hat{A} - M^*\|_F^2 \leq 4k\|A - M^*\|_2^2$ . Hence, the cost of the solution using centers  $\nu_i$ s must be at most  $36k\sigma^2N$  (using 15).

Next, suppose that there exists  $\mu_i$  such that for all  $j$ ,  $\|\mu_i - \nu_j\| > 20\sqrt{k}\alpha \frac{\|A - M^*\|}{\sqrt{Nw_{\min}}}$ . let's compute the cost paid by the points in component  $C_i$  in the clustering obtained via the approximation algorithm. For any  $x \in C_i$  let  $\nu_x$  be the center that it is closest to. Then the cost is at least  $\sum_{x \in C_i} \|x - \nu_x\|^2 \geq \sum_{x \in C_i} \frac{1}{2} \|\mu_i - \nu_x\|^2 - \|x - \mu_i\|^2$ . The first summation is

at least  $\frac{1}{2}|Nw_{\min}|(400\alpha k \frac{\|A - M^*\|^2}{Nw_{\min}}) > 200k\alpha\|A - M^*\|^2$ . The second summation is at most  $\sum_{x \in C_i} \|x - \mu_i\|^2 \leq \sum_i \sum_{x \in C_i} \|x - \mu_i\|^2 = \|\hat{A} - M^*\|_F^2 \leq 4k\|A - M^*\|^2$ . Hence, we reach a contradiction to the fact that the solution obtained via  $\nu_i$ s is an  $\alpha$ -approximation to the optimal cost.  $\square$

The proof of the above theorem already provides a good initializer provided  $\Delta$  is larger than  $\sqrt{k \frac{\log N}{w_{\min}}}$  and one uses a constant factor approximation algorithm for  $k$ -means (Ahmadian et al., 2016). Furthermore, if  $\Delta$  is larger than  $\sqrt{k \log k \frac{\log N}{w_{\min}}}$ , then one can instead use the simpler and faster  $k$ -means++ approximation algorithm (Arthur & Vassilvitskii, 2007). The above lemma has a bad dependence on  $w_{\min}$ . However, using the Boosting technique of (Kumar & Kannan, 2010) we can reduce the dependence to  $\Delta > 25\sqrt{k \log N}$  and hence prove Proposition 8.2. We provide a proof of this in the Appendix.

## 9. Lower Bounds for Semi-random GMMs

We prove the following theorem.

**Theorem 9.1.** *For any  $d, k \in \mathbb{Z}_+$ , there exists  $N_0 = \text{poly}(d, k)$  and a universal constant  $c_1 > 0$  such that the following holds for all  $N \geq N_0$  and  $\Delta$  such that  $\sqrt{\log N} \leq \Delta \leq d/(64 \log d)$ . There exists an instance  $\mathcal{X}$  on  $N$  points in  $d$  dimensions with planted clustering  $C_1, \dots, C_k$  generated by applying semi-random perturbations to points generated from a mixture of spherical Gaussians with means  $\mu_1, \mu_2, \dots, \mu_k$ , covariance  $\sigma^2 I$  and weights being  $1/k$  each, with separation  $\forall i \neq j \in [k], \|\mu_i - \mu_j\|_2 \geq \Delta\sigma$ , such that any locally optimal  $k$ -means clustering solution  $C'_1, C'_2, \dots, C'_k$  of  $\mathcal{X}$  satisfies w.h.p.*

$$\min_{\pi \in \text{Perm}_k} \sum_{i=1}^k |C'_{\pi(i)} \Delta C_i| \geq \frac{c_1 kd}{\Delta^4}.$$

*It suffices to set  $N_0(d, k) := c_0 k^2 d^{3/2} \log^2(kd)$ , where  $c_0 > 0$  is a sufficiently large universal constant.*

**Remark 9.2.** Note that the lower bound also applies in particular to the more general semi-random model in Definition 1.1; in this instance, the points are drawn i.i.d. from the mixture of spherical Gaussians, before applying semi-random perturbations. Further, this lower bound holds for any *locally optimal solution*, and not just the optimal solution.

The lower bound construction will pick an arbitrary  $\Omega(d/\Delta^4)$  points from  $k/2$  clusters, and carefully choose



a semi-random perturbation to all the points so that these  $\Omega(kd/\Delta^4)$  points are misclassified. We start with a simple lemma that shows that an appropriate semi-random perturbation can move the mean of a cluster by an amount  $O(\sigma)$  along any fixed direction.

**Lemma 9.3.** *Consider a spherical Gaussian in  $d$  dimensions with mean  $\mu$  and covariance  $\sigma^2 I$ , and let  $\hat{e}$  be a fixed unit vector. Consider the semi-random perturbation given by*

$$\forall y \in \mathbb{R}^d, h(y) = \begin{cases} \mu & \text{if } \langle y - \mu, \hat{e} \rangle < 0 \\ y & \text{otherwise} \end{cases}.$$

Then we have  $\mathbb{E}[h(y)] = \mu + \frac{1}{\sqrt{2\pi}}\sigma\hat{e}$ .

*Proof.* We assume without loss of generality that  $\mu = 0, \sigma = 1$  (by shifting and scaling) and  $\hat{e} = (1, 0, 0, \dots, 0) \in \mathbb{R}^d$  (by the rotational symmetry of a spherical Gaussian). Let  $\gamma$  be the p.d.f. of the standard Gaussian in  $d$  dimensions with mean 0, and  $\gamma'(y)$  be the distribution on  $y$  conditioned on the event  $[y(1) = \langle y, \hat{e} \rangle > 0]$ . First, we observe that  $\mathbb{E}[h(y)|y_1 < 0] = 0$  from construction, and  $\mathbb{E}[h(y)|y_1 > 0] = \mathbb{E}_{y \sim \gamma'(y)}[y]$ . Further, since the  $(d-1)$  co-ordinates of  $y$  orthogonal to  $\hat{e}$  are independent of  $y_1$ ,

$$\begin{aligned} \mathbb{E}[h(y)] &= \mathbb{P}[y_1 < 0] \mathbb{E}[h(y)|y_1 < 0] \\ &\quad + \mathbb{P}[y_1 > 0] \mathbb{E}[h(y)|y_1 > 0] \\ &= \frac{1}{2} \mathbb{E}[y_1 | y_1 > 0] \hat{e} \\ \mathbb{E}[h(y)] - \mu &= \left( \frac{1}{2\sqrt{2\pi}} \int_{-\infty}^{\infty} |y_1| \exp(-y_1^2/2) dy_1 \right) \hat{e} \\ &= \frac{\sigma}{\sqrt{2\pi}} \hat{e}. \end{aligned}$$

□

**Construction.** Set  $m := c_1 d/\Delta^4$  for some appropriately small constant  $c_1 \in (0, 1)$ . We assume without loss of generality that  $k$  is even (the following construction also works for odd  $k$  by leaving the last cluster unchanged). We pair up the clusters into  $k/2$  pairs  $\{(C_1, C_2), (C_3, C_4), \dots, (C_{k-1}, C_k)\}$ , and we will ensure that  $m$  points are misclassified in each of the  $k/2$  clusters  $C_1, C_3, \dots, C_{k-1}$ . The parameters of the mixture of spherical Gaussians  $\mathcal{G}$  are set up as follows. For each  $i \in \{1, 3, 5, \dots, k-1\}$ ,  $\|\mu_i - \mu_{i+1}\|_2 = \Delta\sigma$ , and all the other inter-mean distances (across different pairs) are at least  $M\sigma$  which is arbitrarily large (think of  $M \mapsto \infty$ ).

- Let for any  $i \in \{1, 3, \dots, k-1\}$ ,  $Z_i \subset C_i$  be the first  $m$  points in cluster  $C_i$  respectively among

the samples  $y^{(1)}, \dots, y^{(N)}$  drawn from  $\mathcal{G}$  (these  $m$  points inside the clusters can be chosen arbitrarily). Set  $Z_i = \emptyset$  for  $i \in \{2, 4, \dots, k\}$ .

- For each  $i \in \{1, 3, \dots, k-1\}$ , set  $\hat{e}_i$  to be the unit vector along  $u_i = \frac{1}{\sigma\sqrt{md}} \sum_{y \in Z_i} (y - \mu_i)$ .
- For each  $i \in \{1, 3, \dots, k-1\}$  apply the following semi-random perturbation given by Lemma 9.3 to points in cluster  $C_{i+1}$  along  $\hat{e}_i$ , i.e., each point  $y^{(t)} \in C_{i+1}$

$$x^{(t)} = h(y^{(t)}) = \begin{cases} \mu_{i+1} & \text{if } \langle y^{(t)} - \mu_{i+1}, \hat{e}_i \rangle < 0 \\ y^{(t)} & \text{otherwise} \end{cases}.$$

Note that the semi-random perturbations are only made to points in the even clusters (based on a few points in its respective odd cluster). The lower bound proof proceeds in two parts. Lemma 9.4 (using Lemma 9.3) and Lemma 9.5 shows that in any  $k$ -means optimal clustering the means of each even cluster  $C_i$  moves by roughly  $\Omega(\sigma) \cdot \hat{e}_{i-1}$ . Lemma 9.6 then shows that these means will classify all the  $m$  points in  $Z_{i-1}$  *incorrectly* w.h.p. In this proof w.h.p. will refer to a probability of at least  $1 - o(1)$  unless specified otherwise (this can be made  $1 - 1/\text{poly}(m, k)$  by choosing suitable constants).

We start with two simple concentration statements about the points in  $Z_i$  (from Lemma A.2 and Lemma 7.6). We have with probability at least  $1 - 1/(mk)$ ,  $\forall i \in \{1, 3, k-1\}$ ,  $\forall t \in Z_i$ ,

$$\|x^{(t)} - \mu_i\|_2 \leq \sigma(\sqrt{d} + 2\sqrt{\log(mk)}) \quad (22)$$

$$|\langle x^{(t)} - \mu_i, \mu_i - \mu_{i+1} \rangle| \leq 2\sqrt{\log(mk)}\Delta\sigma^2 \quad (23)$$

Let  $\tilde{\mu}_1, \dots, \tilde{\mu}_k$  be the (empirical) means of the clusters in the planted clustering  $C_1, C_2, \dots, C_k$  after the semi-random perturbations. The following lemma shows that  $\|\tilde{\mu}_i - \mu_i\|_2 \leq \sigma$ .

**Lemma 9.4.** *There exists a universal constant  $c_3 > 0$  s.t. for the semi-random instance  $\mathcal{X}$  described above, we have that w.h.p.*

$$\forall i \in [k], \tilde{\mu}_i = \begin{cases} \mu_i + \frac{1}{\sqrt{2\pi}}\sigma\hat{e}_{i-1} + z_i & \text{if } i \text{ is even} \\ \mu_i + z_i & \text{if } i \text{ is odd} \end{cases},$$

where  $\|z_i\|_2 \leq c_3\sigma\sqrt{\frac{dk}{N}}$ .

The following lemma shows that if  $C_i, C'_i$  are close, then the empirical means are also close.

**Lemma 9.5.** *Consider any cluster  $C_i$  of the instance  $\mathcal{X}$ , and let  $C'_i$  satisfy  $|C'_i \Delta C_i| \leq m'$ . Suppose  $\tilde{\mu}_i$  and*

$\mu'_i$  are the means of clusters  $C_i$  and  $C'_i$  respectively, then

$$\|\mu'_i - \tilde{\mu}_i\|_2 \leq 4\sigma \cdot \frac{m'}{|C_i|} (\sqrt{d} + 2\sqrt{\log N} + \Delta).$$

The following lemma shows that the Voronoi partition about  $\tilde{\mu}_1, \dots, \tilde{\mu}_k$  (or points close to it) incorrectly classify all points in  $Z_i$  for each  $i \in [k]$ .

**Lemma 9.6.** *Let  $\mu'_1, \mu'_2, \dots, \mu'_k$  satisfy  $\|\mu'_i - \tilde{\mu}_i\|_2 \leq \sigma/(16\sqrt{m}(1 + 2\sqrt{\frac{\log N}{d}}))$ , where  $\tilde{\mu}_i$  is the empirical mean of the points in  $C_i$ . Then, we have w.h.p. that for each  $i \in \{1, 3, \dots, k-1\}$ ,  $\|x - \mu'_i\|_2^2 > \|x - \mu'_{i+1}\|_2^2$ , i.e., every point  $x \in Z_i$  is misclassified.*

*Proof of Theorem 9.1.* Let  $C'_1, \dots, C'_k$  be a locally optimal  $k$ -means clustering of  $\mathcal{X}$ , and suppose  $\sum_i |C'_i \Delta C_i| < mk/2$  (for sake of contradiction). For each  $i \in [k]$ , let  $\tilde{\mu}_i$  be the empirical mean of  $C_i$  and  $\mu'_i$  be the empirical mean of  $C'_i$ . Since  $C'_1, \dots, C'_k$  is a locally optimal clustering, the Voronoi partition given by  $\mu'_1, \dots, \mu'_k$  classifies all the points in agreement with  $C'_1, \dots, C'_k$ .

We will now contradict the local optimality of the clustering  $C'_1, \dots, C'_k$ . Every cluster  $C_i$  has at least  $N/(2k)$  points w.h.p. Hence, for each  $i \in [k]$ , from Lemma 9.5 we have

$$\begin{aligned} \|\mu'_i - \tilde{\mu}_i\|_2 &\leq \sigma(\sqrt{d} + 2\sqrt{\log N} + \Delta) \cdot \frac{4|C_i \Delta C'_i|}{\frac{N}{2k}} \\ &\leq \sigma \cdot \frac{8k^2 m (\sqrt{d} + 2\sqrt{\log N} + \Delta)}{N} \\ &\leq \frac{\sigma}{16\sqrt{m}(1 + \sqrt{(\log N)/d})}. \end{aligned}$$

However, from Lemma 9.6, every point in  $\cup_{i \in [k]} Z_i$  is misclassified by  $\mu'_1, \mu'_2, \dots, \mu'_k$ , i.e., the clustering given the Voronoi partition around  $\mu'_1, \dots, \mu'_k$  differs from  $C_1, \dots, C_k$  on at least  $mk/2$  points in total. But  $\sum_{i \in [k]} |C'_i \Delta C_i| < mk/2$ . Hence, this contradicts the local optimality of the clustering  $C'_1, \dots, C'_k$ .  $\square$

Before we prove Lemma 9.5 and Lemma 9.6, we start with a couple of simple claims about the unit vectors  $\hat{e}_1, \hat{e}_3, \dots, \hat{e}_{k-1}$ .

**Lemma 9.7.** *In the above construction, for every  $i \in \{1, 3, \dots, k-1\}$  we have w.h.p.  $\|\hat{e}_i - u_i\|_2^2 \leq 6\sqrt{m \log(mk)}/d$ . Further, for each  $x \in Z_i$ , we have  $\langle x - \mu_i, \hat{e}_i \rangle \geq \frac{1}{2}\sigma\sqrt{d/m}$ .*

*Proof.* Let us fix an  $i \in \{1, 3, \dots, k-1\}$ . Let  $y^{(1)}, y^{(2)}, \dots, y^{(m)} \in Z_i$  and  $\bar{y}^{(t)} = y^{(t)} - \mu_i$ .

From (22), we know that w.h.p.,  $\|\bar{y}^{(t)}\|_2 \leq \sigma(\sqrt{d} + 2\sqrt{\log m}) \forall t \in [m]$ . Fix  $t \in [m]$ , and let  $Q(t) = \sum_{t' \in [m] \setminus \{t\}} \langle \bar{y}^{(t)}, \bar{y}^{(t')} \rangle$ . For  $t' \neq t$ , due to independence and spherical symmetry,  $\frac{1}{\|\bar{y}^{(t)}\|_2} \langle \bar{y}^{(t)}, \bar{y}^{(t')} \rangle$  is distributed as a normal r.v. with mean 0 and variance  $\sigma^2$ . Further,  $Q(t)/\|y^{(t)}\|_2$  is distributed as a normal r.v. with mean 0 and variance  $\sigma^2 m$ . Hence,

$$\begin{aligned} Q(t) &= \|y^{(t)}\|_2 \cdot \sum_{t' \in [m] \setminus \{t\}} \frac{\langle \bar{y}^{(t')}, \bar{y}^{(t)} \rangle}{\|\bar{y}^{(t)}\|_2} \\ &\leq \sigma^2 (\sqrt{d} + \sqrt{\log(mk)}) \cdot 2\sqrt{m \log(mk)}, \end{aligned} \quad (24)$$

with probability at least  $1 - 1/(mk)^2$ . Hence, w.h.p.  $Q(t) \leq 4\sigma^2 \sqrt{dm \log(mk)}$  for all  $t \in [m]$ .

For the first part, we see that

$$\begin{aligned} \|u_i\|_2^2 &= \frac{1}{\sigma^2 md} \left( \sum_{t \in [m]} \|\bar{y}^{(t)}\|_2^2 + 2 \sum_{t \neq t' \in [m]} \langle \bar{y}^{(t)}, \bar{y}^{(t')} \rangle \right) \\ &= \frac{1}{\sigma^2 md} \left( \sum_{t \in [m]} \|\bar{y}^{(t)}\|_2^2 + 2 \sum_{t \in [m]} Q(t) \right). \end{aligned}$$

Along with (22), the bound on  $Q(t)$  and  $\mathbb{E}[\|y^{(t)}\|_2^2] = d\sigma^2$ , this implies

$$\begin{aligned} \|u_i\|_2^2 - 1 &\leq \frac{1}{md} (4m\sqrt{d \log(mk)} + 4m \log(mk) + 4m\sqrt{dm \log(mk)}) \\ &\quad \text{w.p. at least } 1 - 1/(mk) \end{aligned}$$

$$\|u_i\|_2^2 - 1 \leq 6\sqrt{\frac{m \log(mk)}{d}} \quad \text{with probability at least } 1 - 1/(mk).$$

Since  $\hat{e}_i$  is the unit vector along  $u_i$ , and performing a union bound over all  $i$  we have that w.h.p.,  $\|\hat{e}_i - u_i\|_2^2 \leq 6\sqrt{m \log(mk)}/d$ .

For the furthermore part, suppose  $x = y^{(t)}$  for some  $t \in [m]$  then

$$\begin{aligned} \langle \bar{x}, \hat{e}_i \rangle &= \frac{1}{\|u_i\|_2 \sqrt{md}} \sum_{t' \in [m]} \langle y^{(t)}, y^{(t')} \rangle \\ &\geq \frac{1}{\|u_i\|_2 \sqrt{md}} \left( \|\bar{y}^{(t)}\|_2^2 - \sum_{t' \neq t} |\langle \bar{y}^{(t)}, \bar{y}^{(t')} \rangle| \right) \\ &\geq \frac{\sigma^2}{\|u_i\|_2 \sqrt{md}} \left( (d - \sqrt{d \log(mk)}) - Q(t) \right) \\ &\geq \frac{\sigma^2}{\|u_i\|_2 \sqrt{md}} (d - 4\sqrt{dm \log(mk)}) \\ &\geq \frac{\sigma^2}{4} \sqrt{dm}, \end{aligned}$$

since  $64m \log m \leq d$  and  $\|u_i\|_2 \leq 2$  w.h.p.  $\square$

*Proof of Lemma 9.4.* The lemma follows in a straightforward way from Lemma 9.3 and by standard concentration bounds. Firstly, the clusters  $C_i$  for odd  $i$  are

unaffected by the perturbation. Hence,  $\mathbb{E}_{x \in C_i}[x] = \mu_i$  and from Lemma A.4, the empirical mean of the points in  $C_i$  (there are at least  $N/(2k)$  of them w.h.p.) gives the above lemma. Consider any even  $i$ . From Lemma 9.3, the semi-random perturbation applied to the points in  $C_i$  along the direction  $\widehat{e}_{i-1}$  ensures that  $\mathbb{E}_{x \in C_i}[x] = \mu_i + \frac{\sigma}{\sqrt{2\pi}}\widehat{e}_{i-1}$ . Again by Lemma A.4 applied to the points from  $C_i$ , the lemma follows.  $\square$

*Proof of Lemma 9.5.* Let  $i$  be even (an even cluster). First, we note that from our construction, all the points in  $C'_i \setminus C_i \in C_{i-1}$  w.h.p., since the distance between the means  $\|\mu_i - \mu_j\|_2 \geq M\sigma$  when  $j \notin \{i-1, i\}$ , for  $M$  that is chosen to be appropriately large enough. Further,  $\|\mu_i - \mu_{i-1}\|_2 = \Delta\sigma$ . Let  $\bar{x} = x - \mu_i$  if  $i \in C_i$  and  $\bar{x} = x - \mu_j$  if  $x \in C_j$ . Hence w.h.p.,

$$\forall x \in C'_i \cup C_i, \|x - \mu_i\|_2 \leq \Delta\sigma + \|\bar{x}\|_2 \leq \Delta\sigma + (\sqrt{d} + 2\sqrt{\log N})\sigma.$$

Further,  $\tilde{\mu}_i$  is the empirical mean of all the points in  $C_i$ . Let  $\delta = m'/|C_i|$ .

$$\begin{aligned} \mu'_i - \mu_i &= \frac{\sum_{x \in C_i}(x - \mu_i)}{|C_i|} - \frac{\sum_{x \in C_i \setminus C'_i}(x - \mu_i)}{|C'_i|} \\ &\quad + \frac{\sum_{x \in C'_i \setminus C_i}(x - \mu_i)}{|C'_i|} \\ \mu'_i - \tilde{\mu}_i &= (\mu'_i - \mu_i) - (\tilde{\mu}_i - \mu_i) \\ &= (\mu'_i - \mu_i) + \frac{\sum_{x \in C_i}(x - \mu_i)}{|C_i|} \end{aligned}$$

$$\begin{aligned} \text{Hence, } \mu'_i - \tilde{\mu}_i &= \left(\frac{|C_i|}{|C'_i|} - 1\right)(\tilde{\mu}_i - \mu_i) \\ &\quad - \frac{1}{|C'_i|} \sum_{x \in C_i \setminus C'_i} (x - \mu_i) \\ &\quad + \frac{1}{|C'_i|} \sum_{x \in C'_i \setminus C_i} (x - \mu_i) \end{aligned}$$

$$\begin{aligned} \|\mu'_i - \tilde{\mu}_i\|_2 &\leq \left(\frac{\delta}{1-\delta}\right) \|\tilde{\mu}_i - \mu_i\|_2 \\ &\quad + \left(\frac{2\delta}{1-\delta}\right) \max_{x \in C_i \cup C'_i} \|x - \mu_i\| \\ &\leq \left(\frac{2\delta\sigma}{1-\delta}\right) (1 + \Delta + \sqrt{d} + 2\sqrt{\log N}) \\ &\leq 4\delta\sigma (\Delta + \sqrt{d} + 2\sqrt{\log N}), \end{aligned}$$

where  $\|\mu_i - \tilde{\mu}_i\|_2$  is bounded because of Lemma 9.4. A similar argument follows when  $i$  is odd.  $\square$

*Proof of Lemma 9.6.* Let  $i$  be odd, and consider a

point  $x$  in  $Z_i$ , and let  $\bar{x} = x - \mu_i$ .

$$\begin{aligned} \|x - \mu'_i\|_2^2 - \|x - \mu'_{i+1}\|_2^2 &= \|(x - \mu_i) + \mu_i - \mu'_i\|_2^2 \\ &\quad - \|(x - \mu_{i+1}) + (\mu_{i+1} - \mu'_{i+1})\|_2^2 \\ &= \|x - \mu_i\|_2^2 - \|x - \mu_i + (\mu_i - \mu_{i+1})\|_2^2 \\ &\quad + 2\langle x - \mu_i, \mu_i - \mu'_i \rangle \\ &\quad - 2\langle x - \mu_{i+1}, \mu_{i+1} - \mu'_{i+1} \rangle \\ &\quad + \|\mu_i - \mu'_i\|_2^2 - \|\mu_{i+1} - \mu'_{i+1}\|_2^2 \\ &\geq 2\langle \bar{x}, \mu_i - \mu'_i \rangle + 2\langle \bar{x}, \mu'_{i+1} - \mu_{i+1} \rangle \\ &\quad + 2\langle \bar{x}, \mu_{i+1} - \mu_i \rangle - \Delta^2\sigma^2 \\ &\quad - 2\langle \mu_i - \mu_{i+1}, \mu_{i+1} - \mu'_{i+1} \rangle - \sigma^2 \\ &\geq 2\langle \bar{x}, \mu_i - \mu'_i \rangle + 2\langle \bar{x}, \mu'_{i+1} - \mu_{i+1} \rangle \\ &\quad - 4\Delta\sqrt{\log(mk)}\sigma^2 - \Delta^2\sigma^2 - 2\Delta\sigma^2 - \sigma^2, \end{aligned}$$

where the last inequality follows from (23). From Lemma 9.4, we have

$$\begin{aligned} \mu'_{i+1} - \mu_{i+1} &= (\tilde{\mu}_{i+1} - \mu_{i+1}) + (\mu'_{i+1} - \tilde{\mu}_{i+1}) = \frac{1}{\sqrt{2\pi}}\sigma\widehat{e}_i + z'_{i+1}, \\ \text{where } \|z'_{i+1}\|_2 &\leq \|z_{i+1}\|_2 + \|\mu'_{i+1} - \tilde{\mu}_{i+1}\|_2 \\ &\leq \sigma \cdot \frac{1}{(12\sqrt{m}(1 + 2\sqrt{\log N/d}))}, \end{aligned}$$

since  $N/\sqrt{\log N} \geq Cd^{3/2}km$  for some appropriately large constant  $C > 0$ . Similarly  $\mu'_i - \mu_i = z'_i$ , where  $\|z'_i\|_2 \leq \sigma/(12\sqrt{m}(1 + \sqrt{\log N/d}))$ . Hence, simplifying and applying Lemma 9.7 we get

$$\begin{aligned} \|x - \mu'_i\|_2^2 - \|x - \mu'_{i+1}\|_2^2 &\geq \frac{2}{\sqrt{2\pi}}\langle \bar{x}, \widehat{e} \rangle \sigma - 2|\langle \bar{x}, z'_i \rangle| - 2|\langle \bar{x}, z'_{i+1} \rangle| \\ &\quad - 4\Delta\sqrt{\log(mk)}\sigma^2 - \Delta^2\sigma^2 - \sigma^2 \\ &\geq \sqrt{\frac{d}{2\pi m}} \cdot \sigma^2 - 2\|x\|_2(\|z'_i\|_2 + \|z'_{i+1}\|_2) \\ &\quad - 4\Delta^2\sigma^2 \\ &\geq \sigma^2\sqrt{\frac{d}{2\pi m}} - \sigma^2\sqrt{\frac{d}{9m}} - 4\sigma^2\Delta^2 > 0, \end{aligned}$$

since  $m \leq cd/\Delta^4$  for some appropriate constant  $c$  (say  $c = 16\pi$ ).  $\square$

## 10. Conclusion

In this work we initiated the study of clustering data from a semi-random mixture of Gaussians. We proved that the popular Lloyd's algorithm achieves near optimal error. The robustness of the Lloyd's algorithm for the semi-random model suggests a theoretical justification for its widely documented success in practice. A concrete open question left from our work is to extend our lower bound for locally optimal clusterings to a more general statistical lower bound – this would also imply a separation between recovery guarantees

for the semi-random model and the pure GMM model. Robust analysis under semi-random adversaries for related heuristics such as the EM algorithm and studying semi-random variants for other popular statistical models in machine learning will further improve the gap between our theoretical understanding and observed practical performance of algorithms for such models.

## A. Standard Properties of Gaussians

**Lemma A.1.** *Suppose  $x \in \mathbb{R}$  be generated according to  $N(0, \sigma^2)$ , let  $\Phi(t)$  represent the probability that  $x > t$ , and let  $\Phi^{-1}y$  represent the quantile  $t$  at which  $\Phi(t) \leq y$ . Then*

$$\frac{\frac{t}{\sigma}}{\left(\frac{t^2}{\sigma^2} + 1\right)} e^{-\frac{t^2}{2\sigma^2}} \leq \Phi(t) \leq \frac{\sigma}{t} e^{-\frac{t^2}{2\sigma^2}}. \quad (25)$$

Further, there exists a universal constant  $c \in (1, 4)$  such that

$$\frac{1}{c} \sqrt{\log(1/y)} \leq \frac{t}{\sigma} \leq c \sqrt{\log(1/y)}. \quad (26)$$

Let  $\gamma_d$  be the Gaussian measure associated with a standard Gaussian with mean 0 and variance 1 in each direction. We start with a simple fact about the probability mass of high-dimensional spherical Gaussians being concentrated at around  $\sqrt{d}\sigma$ .

Using concentration bounds for the  $\chi^2$  random variables, we have the following bounds for the lengths of vectors picked according to a standard Gaussian in  $d$  dimensions (see (4.3) in (Laurent & Massart, 2000)).

**Lemma A.2.** *For a standard Gaussian in  $d$  dimensions (mean 0 and variance  $\sigma^2$  in each direction), and any  $t > 0$*

$$\begin{aligned} \mathbb{P}_{x \sim \gamma_d} \left[ \|x\|^2 \geq \sigma^2(d + 2\sqrt{dt} + 2t) \right] &\leq e^{-t}. \\ \mathbb{P}_{x \sim \gamma_d} \left[ \|x\|^2 \leq \sigma^2(d - 2\sqrt{dt}) \right] &\leq e^{-t}. \end{aligned}$$

The following lemma follows from Lemma A.2 and a simple coupling to a spherical Gaussian with variance  $\sigma^2 I$ .

**Lemma A.3.** *Consider any points  $y^{(1)}, \dots, y^{(N)}$  drawn from a Gaussian with mean 0 and variance at most  $\sigma^2$  in each direction. Then with high probability we have*

$$\forall \ell \in [N], \|y^{(\ell)}\|_2 \leq \sigma(\sqrt{d} + 2\sqrt{\log N}).$$

*Proof.* Consider a random vector  $z \in \mathbb{R}^d$  generated from a Gaussian with mean 0 and variance  $\sigma^2$  in each

direction. From Lemma A.2,

$$\begin{aligned} \Pr[\|z\|_2 \geq \sigma(\sqrt{d} + 2\sqrt{\log N})] \\ = \Pr[\|z\|_2^2 \geq \sigma^2(d + 4\sqrt{d \log N} + 4 \log N)] \\ \leq \exp(-2 \log N) < N^{-2}. \end{aligned}$$

Fix  $\ell \in [N]$ . By a simple coupling to the spherical Gaussian random variable  $z$  we have

$$\Pr[\|y^{(\ell)}\|_2 \geq \sigma(\sqrt{d} + 2\sqrt{\log N})] \leq \Pr[\|z\| \geq \sigma(\sqrt{d} + 2\sqrt{\log N})] < N^{-2}.$$

By a union bound over all  $\ell \in [N]$ , the lemma follows.  $\square$

**Lemma A.4** ((Vershynin, 2010), Proposition 5.10). *Let  $Y_i \sim N(\mu, \sigma^2 I_{d \times d})$  for  $i = 1, 2, \dots, N$  where  $N = \Omega\left(\frac{d + \log(\frac{1}{\delta})}{\varepsilon^2}\right)$ . Then, with probability at least  $1 - \delta$  we have that*

$$\left\| \frac{1}{N} \sum_{i=1}^N Y_i - \mu \right\|_2 \leq \sigma \varepsilon.$$

**Lemma A.5** ((Vershynin, 2010), Corollary 5.50). *Let  $Y_i \sim N(\mu, \sigma^2 I_{d \times d})$  for  $i = 1, 2, \dots, N$  where  $N = \Omega\left(\frac{d + \log(\frac{1}{\delta})}{\varepsilon^2}\right)$ . Then, with probability at least  $1 - \delta$  we have that*

$$\left\| \frac{1}{N} \sum_{i=1}^N (Y_i - \mu)(Y_i - \mu)^T - \sigma^2 I \right\| \leq \sigma \varepsilon.$$

## B. Proof of Proposition 8.2

The proof will follow the outline in (Kumar & Kannan, 2010). Given  $N$  points from a semi random mixture  $\mathcal{X}$ , we first randomly partition them into two sets  $S_1$  and  $S_2$  of equal size. Let  $T_1, \dots, T_k$  be the partition induced by the true clustering over  $S_1$  and  $T'_1, \dots, T'_k$  be the partition induced over  $S_2$ . Furthermore, let  $A$  be the  $\frac{N}{2} \times d$  matrix consisting of points in  $S_1$  as rows and  $C$  be the  $\frac{N}{2} \times d$  matrix of the corresponding true centers. It is easy to see that with probability at least  $1 - \delta$ , we will have that

$$\forall r \in [k], \min(|T_r|, |T'_r|) \geq \frac{|C_r|}{4}. \quad (27)$$

Assuming that equations (9) to (17) hold with high probability, we next prove that the following conditions

will also hold with high probability

$$\max_{v: \|v\|=1} \frac{1}{|T_r|} \sum_{x \in T_r} [(x - \mu_r) \cdot v]^2 \leq 4\sigma^2, \forall r \in [k] \quad (28)$$

$$\max_{v: \|v\|=1} \frac{1}{|T'_r|} \sum_{x \in T'_r} [(x - \mu_r) \cdot v]^2 \leq 4\sigma^2, \forall r \in [k] \quad (29)$$

$$\left\| \frac{1}{|T_r|} \sum_{x \in T_r} x - \mu_i \right\| \leq 8\sigma, \forall r \in [k] \quad (30)$$

$$\|A - C\|^2 \leq 4\sigma^2 N \quad (31)$$

To prove (28) notice that  $\frac{1}{|T_r|} \sum_{x \in T_r} [(x - \mu_r) \cdot v]^2 \leq \frac{4}{|C_r|} \sum_{x \in C_r} [(x - \mu_r) \cdot v]^2 \leq 16\sigma^2$  (using 13). Similarly, (29) follows. The proof of (30) follows directly from (14). Finally notice that  $\|A - C\|^2 = \max_{v: \|v\|=1} \sum_r \sum_{x \in T_r} [(x - \mu_r) \cdot v]^2 \leq \max_{v: \|v\|=1} \sum_r \sum_{x \in C_r} [(x - \mu_r) \cdot v]^2 \leq 4\sigma^2 N$  (using (15)).

In the analysis below we will assume that the above equations are satisfied by the random partition. Define a graph  $G = (A \cup B, E)$  where the edge set consists of any pair of points that have a distance of at most  $\gamma = 4\sigma(\sqrt{d} + \sqrt{\log N})$ . Notice that from the definition of  $\gamma$ , any two points from the same true cluster  $C_r$  will be connected by an edge in  $G$  (using 9). Next we map the points in  $A$  to a new  $\frac{N}{2}$  dimensional space as follows. For any row  $A_i$  of  $A$  define  $A'_{i,j} = (A_i - \mu) \cdot (B_j - \mu)$  if  $A_i$  and  $B_j$  are in the same connected component of  $G$ . Otherwise, define  $A'_{i,j} = L$  where  $L$  is a large quantity. Here  $\mu$  denotes the mean of the points in the component in  $G$  to which  $A_i$  belongs to. Let  $\theta_r$  denote the mean of the points in  $T_r$  in the new space. We will show that the new mapping amplifies the mean separation.

**Lemma B.1.** *For all  $r \neq s$ ,  $\|\theta_r - \theta_s\| \geq \Omega(\sqrt{|Nw_{\min}|} k \log N) \sigma^2$ .*

*Proof.* We can assume that points in  $T_r, T'_r$  and  $T_s, T'_s$  belong to the same connected component in  $G$ . Otherwise,  $\|\theta_r - \theta_s\| > L$ . Let  $Q$  be the component to which  $T_r$  and  $T_s$  belong with  $\mu$  being the mean of the points in  $Q$ . Then,  $\|\theta_r - \theta_s\|^2 \geq \sum_{B_j \in Q} [(\mu_r - \mu_s) \cdot (B_j - \mu)]^2$ . Notice that  $(\mu_r - \mu_s) \cdot (\mu_r - \mu_s) = (\mu_r - \mu) \cdot (\mu_r - \mu_s) - (\mu_s - \mu) \cdot (\mu_r - \mu_s)$ . Hence, one of the two terms is at least  $\frac{1}{2} \|\mu_r - \mu_s\|^2$  in magnitude. Without loss of generality assume that  $|(\mu_r - \mu) \cdot (\mu_r - \mu_s)| \geq \frac{1}{2} \|\mu_r - \mu_s\|^2 \geq \frac{125^2}{2} k \log N$ .

Now,  $\|\theta_r - \theta_s\|^2 \geq \sum_{B_j \in T'_r} [(\mu_r - \mu_s) \cdot (B_j - \mu)]^2 = \sum_{B_j \in T'_r} [(\mu_r - \mu_s) \cdot (\mu_r - \mu) - (\mu_r - \mu_s) \cdot (\mu_r - B_j)]^2 \geq \frac{1}{2} |B_j| [(\mu_r - \mu_s) \cdot (\mu_r - \mu)]^2 - \sum_{B_j \in T'_r} [(\mu_r - \mu_s) \cdot (\mu_r - B_j)]^2$ . The first term is at least  $\frac{|T'_r|}{8} \|\mu_r - \mu_s\|^4$  and the second

term (in magnitude) is at most  $4|T'_r| \|\mu_r - \mu_s\|^2 \sigma^2$  (using 29). Substituting the bound on  $\|\mu_r - \mu_s\|$  and using 27, we get that  $\|\theta_r - \theta_s\| = \Omega(\sqrt{|Nw_{\min}|}) (k \log N) \sigma^2$ .  $\square$

Let  $A'$  be the matrix of points in the new space and  $C'$  be the matrix of the corresponding centers. We next bound  $\|A' - C'\|$ .

**Lemma B.2.**  $\|A' - C'\| \leq 24\sigma^2 k (\sqrt{d} + 2\sqrt{\log N}) \sqrt{N}$ .

*Proof.* Let  $Y = A' - C'$ . Then we have that  $\|Y\|^2 \leq \|Y^T Y\| = \max_{v: \|v\|=1} \sum_r \sum_{x \in T_r} [(x - \theta_r) \cdot v]^2$ . Let  $Q_r$  be the connected component in  $G$  that the points in  $T_r$  belong to. Then we can write  $\|Y^T Y\| = \max_{v: \|v\|=1} \sum_r \sum_{x \in T_r} \sum_{B_j \in Q_r} v_j^2 [(x - \mu_r) \cdot (B_j - \mu)]^2 \leq \sum_r \sum_{B_j \in Q_r} v_j^2 \sum_{x \in T_r} [(x - \mu_r) \cdot (B_j - \mu)]^2$ . Using 28, we can bound the inner term as  $\sum_{x \in T_r} [(x - \mu_r) \cdot (B_j - \mu)]^2 \leq 4|T_r| \|B_j - \mu\|^2 \sigma^2$ .

Next notice that because of the way  $G$  is constructed, points within the same connected component have distance at most  $k\gamma$ . Hence,  $\|B_j - \mu\| \leq k\gamma$ . Hence,  $\|Y^T Y\| \leq \sum_r \sum_{B_j \in Q_r} v_j^2 4|T_r| (k^2 \gamma^2) \sigma^2 \leq 4Nk^2 \gamma^2 \sigma^2$ . This gives the desired bound on  $\|Y\| = \|A' - C'\|$ .  $\square$

Combining the previous two lemmas we get that  $\|\theta_r - \theta_s\| \geq \Omega(\sqrt{\frac{|Nw_{\min}|}{d}}) \frac{\|A' - C'\|}{\sqrt{N}}$ . We next run the initialization procedure from Section 8.2 by projecting  $A'$  onto the top  $k$  subspace and running a  $k$ -means approximation algorithm. Let  $\phi_1, \dots, \phi_k$  be the means obtained. Using Lemma 8.6 with  $M^* = C'$ , we get that for all  $r$ ,  $\|\phi_r - \theta_r\| \leq 20\sqrt{k\alpha} \frac{\|A' - C'\|}{\sqrt{|Nw_{\min}|}}$ , where  $\alpha$  is the approximation guarantee of the  $k$ -means approximation used. If  $\Delta > c_0 \sqrt{\min\{k, d\} \log N}$ , then we use a constant factor approximation algorithm (Ahmadian et al., 2016). If  $\Delta > c_0 \sqrt{\min\{k, d\} \log k \log N}$ , then we can use the simpler  $k$ -means++ algorithm (Arthur & Vassilvitskii, 2007)

*Proof of Proposition 8.2.* Assuming  $N = \Omega(\frac{k^2 d^2}{w_{\min}^2})$  we get that for all  $r \neq s$ ,  $\|\phi_r - \phi_s\| \geq 10\sqrt{kd} \frac{\|A' - C'\|}{\sqrt{|Nw_{\min}|}}$ . Let  $P_1, \dots, P_k$  be the clustering of points in  $A'$  obtained by using centers  $\phi_1, \dots, \phi_k$ . Then we have that for each  $r$ ,  $|T_r \Delta P_r| \leq \frac{Nw_{\min}}{10\sqrt{d}}$ , since otherwise the total cost paid by the misclassified points will be more than  $4k\|A' - C'\|^2$ . Next we use the clustering  $P_1, \dots, P_k$  to compute means for the original set of points in  $A$ . Let  $\nu_1, \dots, \nu_k$  be the obtained means. We will show that for all  $r$ ,  $\|\nu_r - \mu_r\| \leq \tau\sigma$ , where  $\tau < \frac{\Delta}{4}$ .

Consider a particular partition  $P_r$  that is uniquely identified with  $T_r$ . Let  $n_{r,r}$  be the number of points that belong to both  $P_r$  and  $T_r$  and  $\mu_{r,r}$  be the mean of

those points. Similarly, let  $n_{r,s}$  be the number of points that belong to  $T_s$  originally but belong to  $P_r$  in the current clustering, and let  $\mu_{r,s}$  be their mean. Then,  $\|\mu_r - \nu_r\| \leq \frac{n_{r,r}}{|P_r|} \|\mu_{r,r} - \mu_r\| + \sum_{s \neq r} \frac{n_{r,s}}{|P_r|} \|\mu_{r,s} - \mu_r\|$ . We can bound  $\|\mu_{r,r} - \mu_r\|$  by  $O(\sigma)$  using 14 and  $\|\mu_{r,s} - \mu_r\|$  by  $O(k(\sqrt{d} + 2\sqrt{\log N}))$  using 9 and the fact that points in  $r$  and  $s$  must belong to the same component in  $G$ . Combining we get the claim.  $\square$

## Acknowledgements

Aravindan Vijayaraghavan is supported by the National Science Foundation (NSF) under Grant No. CCF-1652491 and CCF-1637585.

## References

- Achlioptas, D. and McSherry, F. On spectral learning of mixtures of distributions. In *Learning Theory*, pp. 458–469. Springer, 2005.
- Ackerman, M. and Ben-David, S. Clusterability: A theoretical study. In *Artificial Intelligence and Statistics*, pp. 1–8, 2009.
- Ahmadian, S., Norouzi-Fard, A., Svensson, O., and Ward, J. Better guarantees for k-means and euclidean k-median by primal-dual algorithms. *CoRR*, abs/1612.07925, 2016. URL <http://arxiv.org/abs/1612.07925>.
- Anderson, J., Belkin, M., Goyal, N., Rademacher, L., and Voss, J. R. The more, the merrier: the blessing of dimensionality for learning large Gaussian mixtures. In *Proceedings of The 27th Conference on Learning Theory, COLT 2014, Barcelona, Spain, June 13-15, 2014*, pp. 1135–1164, 2014. URL <http://jmlr.org/proceedings/papers/v35/anderson14.html>.
- Angelidakis, H., Makarychev, K., and Makarychev, Y. Algorithms for stable and perturbation-resilient problems. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 438–451. ACM, 2017.
- Arora, S. and Kannan, R. Learning mixtures of arbitrary Gaussians. In *Proceedings of the thirty-third annual ACM symposium on Theory of computing*, pp. 247–257. ACM, 2001.
- Arthur, D. and Vassilvitskii, S. On the worst case complexity of the k-means method. Technical report, Stanford, 2005.
- Arthur, D. and Vassilvitskii, S. K-means++: The advantages of careful seeding. In *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '07, pp. 1027–1035, 2007. ISBN 978-0-898716-24-5. URL <http://dl.acm.org/citation.cfm?id=1283383.1283494>.
- Awasthi, P. and Sheffet, O. Improved spectral-norm bounds for clustering. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pp. 37–49. Springer, 2012.
- Awasthi, P., Blum, A., and Sheffet, O. Center-based clustering under perturbation stability. *Information Processing Letters*, 112(1-2):49–54, 2012.
- Awasthi, P., Charikar, M., Krishnaswamy, R., and Sinop, A. K. The hardness of approximation of euclidean k-means. *arXiv preprint arXiv:1502.03316*, 2015.
- Balakrishnan, S., Wainwright, M. J., and Yu, B. Statistical guarantees for the EM algorithm: From population to sample-based analysis. *CoRR*, abs/1408.2156, 2014. URL <http://arxiv.org/abs/1408.2156>.
- Balcan, M. F. and Liang, Y. Clustering under perturbation resilience. In *International Colloquium on Automata, Languages, and Programming*, pp. 63–74. Springer, 2012.
- Belkin, M. and Sinha, K. Polynomial learning of distribution families. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pp. 103–112. IEEE, 2010.
- Bhaskara, A., Charikar, M., Moitra, A., and Vijayaraghavan, A. Smoothed analysis of tensor decompositions. In *Proceedings of the 46th Symposium on Theory of Computing (STOC)*. ACM, 2014a.
- Bhaskara, A., Charikar, M., and Vijayaraghavan, A. Uniqueness of tensor decompositions with applications to polynomial identifiability. *Proceedings of the Conference on Learning Theory (COLT)*, 2014b.
- Blum, A. and Spencer, J. Coloring random and semi-random k-colorable graphs. *J. Algorithms*, 19:204–234, September 1995. ISSN 0196-6774. doi: <http://dx.doi.org/10.1006/jagm.1995.1034>. URL <http://dx.doi.org/10.1006/jagm.1995.1034>.
- Brubaker, S. C. Robust PCA and clustering in noisy mixtures. In *Proceedings of the Symposium on Discrete Algorithms*, pp. 1078–1087, 2009.
- Brubaker, S. C. and Vempala, S. Isotropic pca and affine-invariant clustering. In *Proceedings of the 2008 49th Annual IEEE Symposium on Foundations of*

- Computer Science*, FOCS '08, pp. 551–560, Washington, DC, USA, 2008. IEEE Computer Society. ISBN 978-0-7695-3436-7. doi: 10.1109/FOCS.2008.48. URL <http://dx.doi.org/10.1109/FOCS.2008.48>.
- Charikar, M., Steinhardt, J., and Valiant, G. Learning from untrusted data. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2017, pp. 47–60, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4528-6. doi: 10.1145/3055399.3055491. URL <http://doi.acm.org/10.1145/3055399.3055491>.
- Dasgupta, S. Learning mixtures of Gaussians. In *Foundations of Computer Science, 1999. 40th Annual Symposium on*, pp. 634–644. IEEE, 1999.
- Dasgupta, S. and Schulman, L. A probabilistic analysis of EM for mixtures of separated, spherical Gaussians. *The Journal of Machine Learning Research*, 8:203–226, 2007.
- Daskalakis, C., Tzamos, C., and Zampetakis, M. Ten steps of EM suffice for mixtures of two Gaussians. *CoRR*, abs/1609.00368, 2016.
- David, R. and Feige, U. On the effect of randomness on planted 3-coloring models. In *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2016, pp. 77–90, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4132-5. doi: 10.1145/2897518.2897561. URL <http://doi.acm.org/10.1145/2897518.2897561>.
- Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Moitra, A., and Stewart, A. Robust estimators in high dimensions without the computational intractability. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 655–664, Oct 2016. doi: 10.1109/FOCS.2016.85.
- Diakonikolas, I., Kane, D. M., and Stewart, A. List-decodable robust mean estimation and learning mixtures of spherical gaussians. *CoRR*, abs/1711.07211, 2017. URL <http://arxiv.org/abs/1711.07211>.
- Dutta, A., Vijayaraghavan, A., and Wang, A. Clustering stable instances of euclidean k-means. *Proceedings of Neural Information Processing Systems (NIPS)*, 2017.
- Feige, U. and Kilian, J. Heuristics for finding large independent sets, with applications to coloring semi-random graphs. In *Foundations of Computer Science, 1998. Proceedings. 39th Annual Symposium on*, pp. 674–683, nov 1998. doi: 10.1109/SFCS.1998.743518.
- Ge, R., Huang, Q., and Kakade, S. M. Learning mixtures of Gaussians in high dimensions. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pp. 761–770, 2015. doi: 10.1145/2746539.2746616. URL <http://doi.acm.org/10.1145/2746539.2746616>.
- Goyal, N., Vempala, S., and Xiao, Y. Fourier PCA and robust tensor decomposition. In *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pp. 584–593, 2014. doi: 10.1145/2591796.2591875. URL <http://doi.acm.org/10.1145/2591796.2591875>.
- Hopkins, S. B. and Li, J. Mixture models, robustness, and sum of squares proofs. *CoRR*, abs/1711.07454, 2017. URL <http://arxiv.org/abs/1711.07454>.
- Hsu, D. and Kakade, S. M. Learning Gaussian mixture models: Moment methods and spectral decompositions. *arXiv preprint arXiv:1206.5766*, 2012.
- Hsu, D. and Kakade, S. M. Learning mixtures of spherical Gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pp. 11–20. ACM, 2013.
- Kalai, A. T., Moitra, A., and Valiant, G. Efficiently learning mixtures of two Gaussians. In *Proceedings of the 42nd ACM symposium on Theory of computing*, pp. 553–562. ACM, 2010.
- Kannan, R., Salmasian, H., and Vempala, S. The spectral method for general mixture models. *SIAM J. Comput.*, 38(3):1141–1156, 2008. doi: 10.1137/S0097539704445925. URL <http://dx.doi.org/10.1137/S0097539704445925>.
- Kanungo, T., Mount, D. M., Netanyahu, N. S., Piatko, C. D., Silverman, R., and Wu, A. Y. A local search approximation algorithm for k-means clustering. In *Proceedings of the eighteenth annual symposium on Computational geometry*, pp. 10–18. ACM, 2002.
- Kothari, P. K. and Steinhardt, J. Better agnostic clustering via relaxed tensor norms. *CoRR*, abs/1711.07465, 2017. URL <http://arxiv.org/abs/1711.07465>.
- Kumar, A. and Kannan, R. Clustering with spectral norm and the k-means algorithm. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pp. 299–308. IEEE, 2010.

- Lai, K. A., Rao, A. B., and Vempala, S. Agnostic estimation of mean and covariance. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 665–674, Oct 2016. doi: 10.1109/FOCS.2016.76.
- Laurent, B. and Massart, P. Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.*, 28(5):1302–1338, 10 2000. doi: 10.1214/aos/1015957395. URL <http://dx.doi.org/10.1214/aos/1015957395>.
- Lee, E., Schmidt, M., and Wright, J. Improved and simplified inapproximability for k-means. *Information Processing Letters*, 120:40–43, 2017.
- Lloyd, S. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982.
- Makarychev, K., Makarychev, Y., and Vijayaraghavan, A. Approximation algorithms for semi-random partitioning problems. In *Proceedings of the 44th Symposium on Theory of Computing (STOC)*, pp. 367–384. ACM, 2012.
- Makarychev, K., Makarychev, Y., and Vijayaraghavan, A. Sorting noisy data with partial information. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pp. 515–528. ACM, 2013.
- Makarychev, K., Makarychev, Y., and Vijayaraghavan, A. Constant factor approximations for balanced cut in the random pie model. In *Proceedings of the 46th Symposium on Theory of Computing (STOC)*. ACM, 2014.
- Makarychev, K., Makarychev, Y., and Vijayaraghavan, A. Correlation clustering with noisy partial information. *Proceedings of the Conference on Learning Theory (COLT)*, 2015.
- Makarychev, K., Makarychev, Y., and Vijayaraghavan, A. Learning communities in the presence of errors. *Proceedings of the Conference on Learning Theory (COLT)*, 2016.
- Mathieu, C. and Schudy, W. Correlation clustering with noisy input. In *Proceedings of the Twenty-first Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '10, pp. 712–728, Philadelphia, PA, USA, 2010. Society for Industrial and Applied Mathematics. ISBN 978-0-898716-98-6. URL <http://dl.acm.org/citation.cfm?id=1873601.1873659>.
- Moitra, A. and Valiant, G. Settling the polynomial learnability of mixtures of Gaussians. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pp. 93–102. IEEE, 2010.
- Moitra, A., Perry, W., and Wein, A. S. How robust are reconstruction thresholds for community detection. *CoRR*, abs/1511.01473, 2015.
- Ostrovsky, R., Rabani, Y., Schulman, L. J., and Swamy, C. The effectiveness of lloyd-type methods for the k-means problem. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pp. 165–176. IEEE, 2006.
- Pearson, K. Contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London. A*, 185:71–110, 1894.
- Regev, O. and Vijayaraghavan, A. Learning mixtures of well-separated gaussians. In *Proceedings of the 58th Annual IEEE Foundations of Computer Science (FOCS)*. IEEE, 2017.
- Tang, C. and Monteleoni, C. Convergence rate of stochastic k-means. In Singh, A. and Zhu, J. (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 1495–1503, Fort Lauderdale, FL, USA, 20–22 Apr 2017. PMLR. URL <http://proceedings.mlr.press/v54/tang17b.html>.
- Teicher, H. Identifiability of mixtures. *The annals of Mathematical statistics*, 32(1):244–248, 1961.
- Teicher, H. Identifiability of mixtures of product measures. *The Annals of Mathematical Statistics*, 38(4): 1300–1302, 1967.
- Vempala, S. and Wang, G. A spectral algorithm for learning mixture models. *Journal of Computer and System Sciences*, 68(4):841–860, 2004.
- Vershynin, R. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- Williamson, D. P. and Shmoys, D. B. *The Design of Approximation Algorithms*. Cambridge University Press, New York, NY, USA, 1st edition, 2011.
- Xu, J., Hsu, D. J., and Maleki, A. Global analysis of expectation maximization for mixtures of two Gaussians. In *NIPS*, 2016.