

Sharp-SSL: Selective high-dimensional axis-aligned random projections for semi-supervised learning

Tengyao Wang*, Edgar Dobriban[†], Milana Gataric[‡]
and Richard J. Samworth[‡]

*Department of Statistics, London School of Economics

[†]Department of Statistics and Data Science, University of Pennsylvania

[‡]Statistical Laboratory, University of Cambridge

April 19, 2023

Abstract

We propose a new method for high-dimensional semi-supervised learning problems based on the careful aggregation of the results of a low-dimensional procedure applied to many axis-aligned random projections of the data. Our primary goal is to identify important variables for distinguishing between the classes; existing low-dimensional methods can then be applied for final class assignment. Motivated by a generalized Rayleigh quotient, we score projections according to the traces of the estimated whitened between-class covariance matrices on the projected data. This enables us to assign an importance weight to each variable for a given projection, and to select our signal variables by aggregating these weights over high-scoring projections. Our theory shows that the resulting **Sharp-SSL** algorithm is able to recover the signal coordinates with high probability when we aggregate over sufficiently many random projections and when the base procedure estimates the whitened between-class covariance matrix sufficiently well. The Gaussian EM algorithm is a natural choice as a base procedure, and we provide a new analysis of its performance in semi-supervised settings that controls the parameter estimation error in terms of the proportion of labeled data in the sample. Numerical results on both simulated data and a real colon tumor dataset support the excellent empirical performance of the method.

1 Introduction

Semi-supervised learning, where we attempt to assign observations to one of finitely many groups based on partially-labeled training data, represents a core modern statistical challenge. It is sufficiently general to incorporate, at either extreme, the unsupervised case of no labeled training data (clustering) and the supervised setting of fully-labeled training data (classification). Such tasks abound in many application areas, including genomics (e.g., Eisen et al., 1998), image processing (Jain and Flynn, 1996; Cheplygina, de Bruijne and Pluim, 2019), natural language processing (Liang, 2005; Turian, Ratinov and Bengio, 2010) and anomaly detection (Akçay, Atapour-Abarghouei and Breckon, 2019; Wang et al., 2019). Entry points to the literature on semi-supervised learning

include Zhu (2005), Zhu and Goldberg (2009), Chapelle, Schölkopf and Zien (2006) and Van Engelen and Hoos (2020). For introductions to clustering, see Xu and Wunsch (2005), Kaufman and Rousseeuw (2009) and Xu and Tian (2015), and for classification, see Devroye, Györfi and Lugosi (2013) and Hastie, Tibshirani and Friedman (2009).

A common feature of contemporary semi-supervised learning problems is high-dimensionality, since we may record many covariates having a possible association with the labels corresponding to different observations. This represents a significant challenge, as can be seen by considering a simple two-class problem with more covariates than observations. For any given assignment of class labels, if no subset of n_0 observations lies in an $(n_0 - 2)$ -dimensional affine space, then we can find hyperplanes with orthogonal normal vectors, each of which achieves zero training error (in other words, they perfectly separate the classes). Nevertheless, even in the simple setting where the true Bayes decision boundary is linear, many such hyperplanes may be little better than a random guess on test data.

An appealing approach to tackling high-dimensionality is via random projections into lower-dimensional spaces. Such projections may almost preserve the pairwise distances between observations, as seen from the Johnson–Lindenstrauss lemma (Johnson and Lindenstrauss, 1984; Dasgupta and Gupta, 2003). Moreover, in cases where we have reason to believe that only a relatively small proportion of the variables recorded are relevant for the learning task, we can choose our random projections to be axis-aligned in order to preserve this structure. A third benefit is the possibility of aggregating results over multiple random projections, though this must be done with care so as to avoid noise accumulation. These attractions have meant that random projections have now been employed in many high-dimensional statistical problems, including precision matrix estimation (Marzetta, Tucci and Simon, 2011), two-sample mean testing (Lopes, Jacob and Wainwright, 2011), classification (Durrant and Kabán, 2015; Cannings and Samworth, 2017), (sparse) principal component analysis (Yang et al., 2021; Gataric, Wang and Samworth, 2020), linear regression (Thanei, Heinze and Meinshausen, 2017; Slawski, 2018; Dobriban and Liu, 2019; Ahfock, Astle and Richardson, 2021), clustering (Dasgupta, 1999; Fern and Brodley, 2003; Han and Boutin, 2015; Yellamraju and Boutin, 2018; Anderlucci, Fortunato and Montanari, 2022) and dimensionality reduction (Bingham and Mannila, 2001; Reeve, Kabán and Bootkrajang, 2022). See Cannings (2021) for a review of recent developments in the area.

In this paper, we propose a new method, called **Sharp-SSL** (short for **S**elective **h**igh-dimensional, **a**xis-aligned **r**andom **p**rojections for **S**emi-**S**upervised **L**earning). Our primary goal is to identify a small subset of variables that are particularly helpful for label assignment; existing low-dimensional methods can then be used to complete the learning task. To this end, we generate a large number of axis-aligned random projections, and apply a base learning procedure such as a semi-supervised version of the Gaussian Expectation–Maximization (EM) algorithm to our projected data. Motivated by the notion of a generalized Rayleigh quotient (see (2) below for a formal definition), and to avoid the noise accumulation issue mentioned above, we score the projections by computing the trace of the corresponding estimated whitened between-class covariance matrices. This enables us to assign an importance weight to each variable for a given projection, and we select our signal variables by aggregating these importance weights over the high-scoring projections. See Section 2 for a more detailed description of our methodology.

Section 3 is devoted to a theoretical analysis of our **Sharp-SSL** algorithm. We first show in Theorem 2 that provided the low-dimensional base learning procedure satisfies

a guarantee on the proximity of the estimated whitened between-class covariance matrix to its population analogue, the corresponding high-dimensional semi-supervised learning algorithm can recover the signal coordinates with high probability when we aggregate over sufficiently many random projections. It turns out that both Linear Discriminant Analysis and an EM algorithm are examples of low-dimensional learning procedures that satisfy this proximity guarantee, as we prove in Theorems 3 and 6 respectively. The latter is particularly challenging, and one of the main novel contributions of our analysis is to provide a guarantee on the performance of a d -dimensional Gaussian EM algorithm in a semi-supervised setting. In particular, we control the parameter estimation error in terms of the proportion of labeled data in the sample, showing that with a sample size of n it smoothly interpolates between the $(d/n)^{1/4}$ rate for unsupervised learning and the $(d/n)^{1/2}$ rate for fully-labeled data, up to logarithmic factors. An advantage of the modular approach to our analysis is that it illustrates the way in which the **Sharp-SSL** algorithm can be combined with different base learning algorithms to adapt to different problem settings and reflect the preferences of the practitioner.

In Section 4, we study the numerical performance of the **Sharp-SSL** algorithm. Our first goal, in Section 4.1, is to study the effect of the choices of input parameters to our method, which allows us to recommend sensible default choices for application in our subsequent comparisons. Section 4.2 presents the results of a simulation study involving the **Sharp-SSL** method, as well as five alternative approaches, on high-dimensional clustering tasks (since not all of the competing methods are able to leverage partial label information). We find that the **Sharp-SSL** algorithm is able to attain a misclustering rate very close to that of the optimal Bayes classifier, even with only around 50 observations per cluster, in settings where these alternative techniques may perform poorly. In Section 4.3, we investigate the extent to which the different versions of the **Sharp-SSL** method are able to leverage partial label information. The results here are consistent with the phase transition phenomenon articulated by our theory. Finally, in Section 4.4, we apply the **Sharp-SSL** algorithm, as well as the other methods from our simulation study, on a colon tumor dataset, where we withhold the true labels from the algorithms in order to assess performance. Our analysis supports the ability of the **Sharp-SSL** algorithm to identify signal coordinates (genes) that are useful for identifying patients with and without tumors.

In the broader literature on high-dimensional learning problems, a large number of methods have been developed to leverage sparse low-dimensional structures for both clustering (Witten and Tibshirani, 2010; Azizyan, Singh and Wasserman, 2013; Wasserman, Azizyan and Singh, 2014; Azizyan, Singh and Wasserman, 2015; Jin and Wang, 2016; Verzelen and Arias-Castro, 2017; Löffler, Wein and Bandeira, 2022; Löffler, Zhang and Zhou, 2021) and classification (Cai and Liu, 2011; Witten and Tibshirani, 2011; Mai, Zou and Yuan, 2012; Cai and Zhang, 2019). These methods are not designed for partially-labeled (semi-supervised) settings. Another common approach is to project the data into the span of the top few principal components, and run a standard low-dimensional method such as k -means clustering or the EM algorithm (Butler et al., 2018). This approach can fail if the directions of largest variation in the data are not aligned with the directions separating the clusters. Finally, recent developments in other aspects of semi-supervised learning include self-training (Oymak and Gulcu, 2020), mean estimation (Zhang, Brown and Cai, 2019), choice of k in k -nearest neighbour classification (Cannings, Berrett and Samworth, 2020) and linear regression (Chakraborty and Cai, 2018).

Proofs of all of our results are provided in Section 5, and we conclude this introduction

with some notation used throughout the paper. We write $\mathbb{S}^{d \times d}$ for the set of d -dimensional symmetric matrices, write $\mathbb{S}_+^{d \times d}$ for the subset that are invertible, and write $\mathbb{S}_{K-1}^{d \times d}$ the subset of matrices in $\mathbb{S}^{d \times d}$ of rank at most $K-1$. We write $\mathbb{R}^{d \times d}$ for the set of d -dimensional matrices. For $p \geq d$, let $\mathbb{O}^{p \times d}$ denote the set of $p \times d$ matrices with orthonormal columns. The Euclidean norm is denoted by $\|\cdot\|$, and the operator norm of a matrix is denoted by $\|\cdot\|_{\text{op}}$, so that $\|A\|_{\text{op}} := \sup_{\{x: \|x\|=1\}} \|Ax\|$. Given two sequences (a_n) and (b_n) , we write $a_n \lesssim b_n$ when there exists a universal constant $C > 0$ such that $a_n \leq Cb_n$, and, given an additional problem parameter R , we write $a_n \lesssim_R b_n$ when there exists $C > 0$, depending only on R , such that $a_n \leq Cb_n$.

For any set $S \subseteq \mathbb{R}^d$ and $d \leq |S|$, we write $\binom{S}{d} := \{A \subseteq S : |A| = d\}$. If $S \subseteq \mathbb{R}^d$, we define $\text{sargmax } S$ to be the smallest element in the argmax in the lexicographic order. For a positive integer k , we define $[k] := \{1, \dots, k\}$. For a vector $v = (v_1, \dots, v_k)^\top \in \mathbb{R}^k$, and $j \in [k]$, we define $v_{-j} = (v_1, \dots, v_{j-1}, v_{j+1}, \dots, v_k)^\top \in \mathbb{R}^{k-1}$.

2 The Sharp-SSL algorithm

In this section, we describe in detail the **Sharp-SSL** algorithm for K -class semi-supervised learning, with $K \geq 2$. We aim to provide a unified treatment of clustering, semi-supervised learning and classification. To this end, we assume that for $i \in [n]$, the observation $x_i \in \mathbb{R}^p$ has a true label $y_i^* \in [K]$, but it may be the case that we do not observe y_i^* . Instead, we assume that our observed label y_i takes values in $[K] \cup \{0\}$, where $y_i := y_i^*$ when the true class label is observed, and $y_i := 0$ otherwise. Thus, our data can be regarded as $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^p \times ([K] \cup \{0\})$, and our goal is to construct a *data-dependent classifier*¹, i.e. a Borel measurable function $C : \mathbb{R}^p \times (\mathbb{R}^p \times ([K] \cup \{0\}))^n \rightarrow [K]$, with the interpretation that $C(x; (x_1, y_1), \dots, (x_n, y_n))$ is the predicted class of $x \in \mathbb{R}^p$.

To motivate our **Sharp-SSL** algorithm, it is instructive first to consider a canonical Gaussian classification problem, where our data can be regarded as n independent realizations of a pair (X, Y) taking values in $\mathbb{R}^p \times [K]$, with prior probability $\pi_k := \mathbb{P}(Y = k)$ for the k th class and $X | Y = k \sim \mathcal{N}_p(\nu_k, \Sigma_w)$, for class means $\nu_1, \dots, \nu_K \in \mathbb{R}^p$ and *within-class covariance matrix* $\Sigma_w \in \mathbb{S}_+^{p \times p}$. Let $\nu := \sum_{k=1}^K \pi_k \nu_k \in \mathbb{R}^p$ denote the grand population mean, let

$$\Sigma_b := \sum_{k=1}^K \pi_k (\nu_k - \nu)(\nu_k - \nu)^\top \in \mathbb{S}_{K-1}^{p \times p} \quad (1)$$

denote the *between-class covariance matrix*, and consider $D \in \mathbb{O}^{p \times (K-1)}$ with a column space spanned by $\Sigma_w^{-1}(\nu_1 - \nu), \dots, \Sigma_w^{-1}(\nu_K - \nu)$. Observe that for $k \neq \ell$,

$$\log \left\{ \frac{\mathbb{P}(Y = k | X = x)}{\mathbb{P}(Y = \ell | X = x)} \right\} = \log \left(\frac{\pi_k}{\pi_\ell} \right) - \frac{1}{2} (\nu_k + \nu_\ell)^\top \Sigma_w^{-1} (\nu_k - \nu_\ell) + x^\top \Sigma_w^{-1} (\nu_k - \nu_\ell),$$

from which we deduce that this likelihood ratio, and hence the Bayes classifier $x \mapsto \text{argmax}_{k \in [K]} \mathbb{P}(Y = k | X = x)$, only depends on x through $D^\top x$. Thus, for the purposes of classification, no signal would be lost (and the noise would be reduced) if X were replaced with $D^\top X$.

¹It is convenient to use the term ‘classifier’ here, even though some or all of the labels may be unobserved.

In high-dimensional settings with $p \gg n$, the matrix Σ_w^{-1} is not consistently estimable in general, but we can nevertheless make progress if the vectors $\Sigma_w^{-1}(\nu_1 - \nu), \dots, \Sigma_w^{-1}(\nu_K - \nu)$ are sparse. In other words, writing S_0 for the union of the set of coordinates for which these vectors are non-zero, we suppose that $|S_0| \ll p$; this is a very common assumption in high-dimensional LDA (e.g. Cai and Liu, 2011; Witten and Tibshirani, 2011; Mai, Zou and Yuan, 2012; Cai and Zhang, 2019).

In such a setting, the column space of D has a sparse basis, so it is natural to consider projecting the data onto a small subset of its coordinates. For $d \in [p]$, define the set of axis-aligned projection matrices $\mathcal{P}_d := \{P \in \{0, 1\}^{d \times p} : PP^\top = I_d\}$, i.e. the set of binary $d \times p$ matrices with orthonormal rows. By the argument above, if $d \geq |S_0|$ then there exists $P^* \in \mathcal{P}_d$ such that the error of the Bayes classifier is unchanged by projecting the data along P^* . In practice, it would typically be computationally too expensive to enumerate through all $p(p-1) \cdots (p-d+1)$ axis-aligned projections. Instead, we consider a randomly chosen subset of projections within \mathcal{P}_d . An axis-aligned projection chosen uniformly at random is unlikely to capture all the signal coordinates S_0 , but by aggregating over a carefully-chosen subset of these random projections, we can nevertheless recover the set of signal coordinates under suitable conditions; see Theorem 2 below. To describe our method for choosing good projections, for $V \in \mathbb{O}^{p \times d}$, we define the *generalized Rayleigh quotient* along V by

$$J(V; \Sigma_b, \Sigma_w) := \text{tr}\{(V^\top \Sigma_w V)^{-1}(V^\top \Sigma_b V)\}. \quad (2)$$

Proposition 1 below motivates seeking to choose projections to maximize the generalized Rayleigh quotient by showing that the column span of any maximizer $J(V; \Sigma_b, \Sigma_w)$ over $V \in \mathbb{O}^{p \times d}$ must contain the column space of D .

Proposition 1. *Let $K \geq 2$ and $d \geq K - 1$. Assume that the convex hull of ν_1, \dots, ν_K is $(K - 1)$ -dimensional, and let $V^* \in \text{argmax}_{V \in \mathbb{O}^{p \times d}} J(V; \Sigma_b, \Sigma_w)$. Then the column space of V^* contains the eigenspace corresponding to the $K - 1$ non-zero eigenvalues² of $\Sigma_w^{-1} \Sigma_b$, which is equal to the space spanned by $(\Sigma_w^{-1}(\nu_k - \nu) : k \in [K])$.*

Based on Proposition 1, a natural conceptual approach to maximizing the generalized Rayleigh quotient is to compute the leading $(K - 1)$ -dimensional eigenspace of $\Sigma_w^{-1} \Sigma_b$. This strategy, however, runs into difficulties when we replace these population quantities with their sample versions in the setting of the opening paragraph of this section. More precisely, writing $n_k := \sum_{i=1}^n \mathbb{1}_{\{y_i=k\}}$ for $k \in [K]$, as well as

$$\begin{aligned} \tilde{\Sigma}_w &:= \frac{1}{n} \sum_{k=1}^K \sum_{i=1}^n (x_i - \hat{\nu}_k)(x_i - \hat{\nu}_k)^\top \mathbb{1}_{\{y_i=k\}} \in \mathbb{R}^{p \times p} \\ \tilde{\Sigma}_b &:= \sum_{k=1}^K \frac{n_k}{n} (\hat{\nu}_k - \hat{\nu})(\hat{\nu}_k - \hat{\nu})^\top \in \mathbb{R}^{p \times p}, \end{aligned}$$

for the sample versions of the within-class and between-class covariance matrices respectively, the matrix $\tilde{\Sigma}_w$ is not invertible whenever $p > n$. Fortunately, though, this issue can be resolved by working with the projected data, as long as we choose $d \leq n - K$: the projected data $\{PX_i : i \in [n]\}$ has within-class covariance matrix $P\tilde{\Sigma}_w P^\top \in \mathbb{R}^{d \times d}$

²Even though $\Sigma_w^{-1} \Sigma_b$ is not guaranteed to be symmetric, it is similar (i.e. conjugate) to the symmetric matrix $\Sigma_w^{-1/2} \Sigma_b \Sigma_w^{-1/2}$, so has real eigenvalues and eigenvectors.

and between-class covariance matrix $P\Sigma_b P^\top \in \mathbb{R}^{d \times d}$, so with probability one, the sample version $P\tilde{\Sigma}_w P^\top$ is invertible.

Returning to the general setting of the opening paragraph of this section, then, we seek projections P with large $J(P^\top; \tilde{\Sigma}_b, \tilde{\Sigma}_w) = \text{tr}\{(P\tilde{\Sigma}_w P^\top)^{-1}(P\tilde{\Sigma}_b P^\top)\}$. To this end, for fixed $A, B \in \mathbb{N}$, we sample a set of projections $\{P^{a,b} : a \in [A], b \in [B]\}$ uniformly at random from \mathcal{P}_d . For each a and b , we apply a low-dimensional base algorithm $\psi : (\mathbb{R}^d \times ([K] \cup \{0\}))^n \rightarrow \mathbb{S}_+^{d \times d}$ to the projected data $(P^{a,b}x_1, y_1), \dots, (P^{a,b}x_n, y_n)$ to obtain an estimator $\hat{Q}^{a,b}$ of $(P^{a,b}\Sigma_w P^{a,b,\top})^{-1}(P^{a,b}\Sigma_b P^{a,b,\top})$, the *whitened between-class covariance matrix* of the projected data. We assume throughout for convenience that ψ is *permutation equivariant* in the sense that $\psi((\Pi z_1, y_1), \dots, (\Pi z_n, y_n)) = \Pi\psi((z_1, y_1), \dots, (z_n, y_n))\Pi^\top$ for every permutation matrix $\Pi \in \mathbb{R}^{d \times d}$. One choice for the base algorithm is to set $\hat{Q}^{a,b} = (\hat{\Sigma}_w^{a,b})^{-1}\hat{\Sigma}_b^{a,b}$, where $\hat{\Sigma}_w^{a,b}$ and $\hat{\Sigma}_b^{a,b}$ are estimated projected within- and between-class (or cluster) covariance matrices.

To select projections, for each $a \in [A]$, we define

$$b^*(a) := \underset{b \in [B]}{\text{sargmax}} \text{tr}(\hat{Q}^{a,b}),$$

and select $P^{a,b^*(a)}$. The main rationale for dividing the projections into A groups and selecting one within each group—as opposed to selecting the A projections with the largest values of $\text{tr}(\hat{Q}^{a,b})$ —is that, conditional on the original data, the selected projections are independent and identically distributed. This facilitates our theoretical analysis by enabling the application of concentration inequalities in the proof of Theorem 2.

The diagonal entries of $\{\hat{Q}^{a,b^*(a)} : a \in [A]\}$ measure the importance of the projected variables for the semi-supervised learning task. These can then be converted into importance scores for the original variables by ‘back-projecting’ into the higher-dimensional space, i.e. by forming the vector $\hat{w} = (\hat{w}_1, \dots, \hat{w}_p)^\top \in \mathbb{R}^p$ given by

$$\hat{w}_j := \frac{1}{A} \sum_{a=1}^A \left[P^{a,b^*(a),\top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)} \right]_{j,j}, \quad j \in [p].$$

Finally, we rank the variables by their importance scores, and our estimate \hat{S} of the set of signal coordinates is given by the largest ℓ entries in \hat{w} , breaking ties arbitrarily if necessary, where $\ell \in [p]$ is specified by the practitioner. Pseudocode for the **Sharp-SSL** procedure is given in Algorithm 1.

After applying Algorithm 1 to obtain an estimated set \hat{S} of signal variables, we can then apply any existing semi-supervised learning method for low-dimensional data with input $(P_{\hat{S}}x_i, y_i)_{i \in [n]}$, where $P_{\hat{S}}$ denotes the projection onto the coordinates in \hat{S} .

2.1 Base learning methods

Algorithm 1 relies on a base learning method for low-dimensional data to estimate the projected whitened between-class covariance matrix from the projected data. When all or almost all of the input data are labeled, we can use the procedure outlined in Algorithm 2, which ignores any unlabeled data, for this purpose. On the other hand, when we have a substantial amount of unlabeled data, Algorithm 2 may be inaccurate. In such circumstances, it may be preferable to use Algorithm 3, which runs an *Expectation–Maximization* (EM) procedure to predict the unobserved labels and subsequently estimate the whitened between-class covariance matrix. More precisely, from M random initializations of the

Algorithm 1: Sharp-SSL: Clustering and semi-supervised learning via ensembles of axis-aligned random projections.

Input: Data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^p \times ([K] \cup \{0\})$ (where $y_i = 0$ denotes a missing label);
 Projected dimension $d \in [\min(p, n - K)]$, number of selected signal coordinates $\ell \in [p]$;
 Number $A \in \mathbb{N}$ of groups of projections, number $B \in \mathbb{N}$ of projections in each group;
 Permutation equivariant base algorithm $\psi : (\mathbb{R}^d \times ([K] \cup \{0\}))^n \rightarrow \mathbb{S}_+^{d \times d}$.

Generate axis-aligned random projections $\{P^{a,b} : a \in [A], b \in [B]\}$ independently and uniformly from \mathcal{P}_d .

```

for  $a \in [A]$  do
  | for  $b \in [B]$  do
  | | Let  $\hat{Q}^{a,b} := \psi((P^{a,b}x_1, y_1), \dots, (P^{a,b}x_n, y_n))$ .
  | end
  | Set  $b^*(a) := \operatorname{sargmax}_{b \in [B]} \operatorname{tr}(\hat{Q}^{a,b})$ .
end

```

Let $\hat{w} = (\hat{w}_1, \dots, \hat{w}_p)^\top$, where $\hat{w}_j := \frac{1}{A} \sum_{a=1}^A [P^{a,b^*(a)}, \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j,j}$.

Output: $\hat{S} \subseteq [p]$, defined as the index set of the ℓ largest components of \hat{w} , breaking ties randomly.

cluster means and the within-class covariance matrix, Algorithm 3 uses the EM algorithm to update these quantities, and thereby compute the whitened between-cluster sample covariance matrix estimators $\{\hat{Q}^{[m]} = (\hat{\Sigma}_w^{[m]})^{-1}(\hat{\Sigma}_b^{[m]}) : m \in [M]\}$. We select $\hat{m} \in [M]$ such that $\hat{Q}^{[\hat{m}]}$ is in best agreement with results from the other EM runs; this is made precise in (6).

The algorithm also allows the practitioner to incorporate prior knowledge about the true cluster means and within-cluster covariance matrices, both through optimizing over a restricted constraint set \mathcal{C} in the M step of the EM algorithm, and through the choice of a distribution supported on \mathcal{C} for the initialization of these quantities. An alternative to the EM algorithm for unsupervised learning would be to apply k -means clustering as a base procedure. Previous studies have suggested that these approaches have comparable empirical performance (e.g., de Souto et al., 2008; Rodriguez et al., 2019, and references therein), but the EM algorithm is more amenable to theoretical analysis in our setting.

3 Theoretical guarantees

3.1 Results for the high-level algorithm

In this subsection, we consider independent triples $(X_1, Y_1, Y_1^*), \dots, (X_n, Y_n, Y_n^*)$ taking values in $\mathbb{R}^p \times ([K] \cup \{0\}) \times [K]$. We recall that Y_i^* denotes the true label of the i th observation, and that $Y_i := Y_i^*$ if the i th label is observed, and $Y_i := 0$ otherwise. For $k \in [K]$, let $\pi_k := \mathbb{P}(Y_1^* = k)$ and $\nu_k^* := \mathbb{E}(X_1 | Y_1^* = k)$ denote the prior probability and the cluster mean of the k th cluster respectively, let $\nu^* := \sum_{k=1}^K \pi_k \nu_k^*$ denote the weighted cluster mean and let $\Sigma_w := \operatorname{Cov}(X_1 | Y_1^* = k)$ denote the common within-cluster covariance matrix. With the between-cluster covariance matrix Σ_b from (1), our

Algorithm 2: Base learning using only labeled data

Input: $(z_1, y_1), \dots, (z_n, y_n) \in \mathbb{R}^d \times ([K] \cup \{0\})$, where the convex hull of $(z_i)_{i:y_i=k}$ is d -dimensional for at least one $k \in [K]$.

Compute $\hat{\mu} := n^{-1} \sum_{i=1}^n z_i$

for $k \in [K]$ **do**

 Set $n_k := |\{i : y_i = k\}|$ and $n' := \sum_{k=1}^K n_k$

 Compute $\hat{\mu}_k := n_k^{-1} \sum_{i:y_i=k} z_i$ (with the convention that $\hat{\mu}_k := 0$ if $n_k = 0$).

end

Compute the within-class and between-class covariance matrices as

$$\hat{\Sigma}_w := \frac{1}{n'} \sum_{i=1}^n (z_i - \hat{\mu}_{y_i})(z_i - \hat{\mu}_{y_i})^\top \quad \text{and} \quad \hat{\Sigma}_b := \sum_{k=1}^K \frac{n_k}{n'} (\hat{\mu}_k - \hat{\mu})(\hat{\mu}_k - \hat{\mu})^\top. \quad (3)$$

Output: $\hat{Q} := \hat{\Sigma}_w^{-1} \hat{\Sigma}_b$.

goal is to estimate the set of signal coordinates,

$$S_0 := \{j \in [p] : (\Sigma_w^{-1} \Sigma_b)_{j,j} \neq 0\},$$

and we write $s_0 := |S_0|$.

Our first main theoretical result shows that if the base algorithm is accurate on each low-dimensional projection and A is large, then with high probability, all signal coordinates are selected.

Theorem 2. Define $\gamma_{\min} := \min_{j \in S_0} (\Sigma_w^{-1} \Sigma_b)_{j,j}$ and $\gamma_{\max} := \max_{j \in S_0} (\Sigma_w^{-1} \Sigma_b)_{j,j}$. Let \hat{S} be the output of Algorithm 1 with input $K, p, (X_1, Y_1), \dots, (X_n, Y_n), A, B, d \geq s_0, \ell \geq s_0$ and permutation equivariant base procedure ψ . Write

$$\varepsilon := \mathbb{P} \left(\max_{P \in \mathcal{P}_d} \left\| \psi((PX_i, Y_i)_{i \in [n]}) - P \Sigma_w^{-1} \Sigma_b P^\top \right\|_{\text{op}} \geq \frac{\gamma_{\min}}{4(K-1)} \right). \quad (7)$$

Then

$$\mathbb{P}(S_0 \subseteq \hat{S}) \geq 1 - \varepsilon - p e^{-A \gamma_{\min}^2 / (50 p^2 \gamma_{\max}^2)}.$$

In fact, we can see from the proof of Theorem 2 that the following stronger conclusion holds: for any realization $(x_i, y_i)_{i \in [n]}$ of the data satisfying

$$\max_{P \in \mathcal{P}_d} \left\| \psi((Px_i, y_i)_{i \in [n]}) - P \Sigma_w^{-1} \Sigma_b P^\top \right\|_{\text{op}} < \frac{\gamma_{\min}}{4(K-1)}, \quad (8)$$

we have $\mathbb{P}(S_0 \subseteq \hat{S} \mid (X_i, Y_i)_{i \in [n]} = (x_i, y_i)_{i \in [n]}) \geq 1 - p e^{-A \gamma_{\min}^2 / (50 p^2 \gamma_{\max}^2)}$. Note here that, after conditioning on the data, the probability is taken over the randomness in the projections. An attraction of Theorem 2 is its generality, and in particular the fact that we do not impose strong distributional assumptions — we simply require control of ε in (7). The price we pay for this generality is that the probability bound may be loose in particular cases; for example, the bound holds even with $B = 1$, though in practice we would expect it to improve as B increases.

Algorithm 3: Base learning using partially labeled data via an EM algorithm

Input: Data $(z_1, y_1), \dots, (z_n, y_n) \in \mathbb{R}^d \times ([K] \cup \{0\})$. A constraint set $\mathcal{C} \subseteq (\mathbb{R}^d)^K \times \mathbb{S}_+^{d \times d}$ and a probability distribution $\pi_{\mathcal{C}}$ supported on \mathcal{C} .
 Number of random initializations M . Number of iterations T .

for $m \in [M]$ **do**

Randomly sample $(\hat{\mu}_1, \dots, \hat{\mu}_K, \hat{\Sigma}_w) \sim \pi_{\mathcal{C}}$.

for $t \in [T]$ **do**

(E step) Compute the soft-label matrix $(L_{i,k})_{i \in [n], k \in [K]}$

$$L_{i,k} := \left(\frac{e^{-\frac{1}{2}(z_i - \hat{\mu}_k)^\top \hat{\Sigma}_w^{-1} (z_i - \hat{\mu}_k)}}{\sum_{\ell=1}^K e^{-\frac{1}{2}(z_i - \hat{\mu}_\ell)^\top \hat{\Sigma}_w^{-1} (z_i - \hat{\mu}_\ell)}} \right) \mathbb{1}_{\{y_i=0\}} + \mathbb{1}_{\{y_i=k\}}. \quad (4)$$

(M step) Update parameter estimates by

$$\begin{aligned} & (\hat{\mu}_1, \dots, \hat{\mu}_K, \hat{\Sigma}_w) \\ & := \operatorname{argmin}_{(\mu_1, \dots, \mu_K, \Sigma) \in \mathcal{C}} \left\{ \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K L_{i,k} (z_i - \mu_k)^\top \Sigma^{-1} (z_i - \mu_k) + \log \det \Sigma \right\}. \end{aligned} \quad (5)$$

end

Compute $(L_{i,k})_{i \in [n], k \in [K]}$ using the final values of $(\hat{\mu}_1, \dots, \hat{\mu}_K, \hat{\Sigma}_w)$ as in (4).

Compute $\hat{\mu}_{\text{tot}} := n^{-1} \sum_{i=1}^n \sum_{k=1}^K L_{i,k} \hat{\mu}_k$ and the between-class covariance matrix

$$\hat{\Sigma}_b := \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K L_{i,k} (\hat{\mu}_k - \hat{\mu}_{\text{tot}})(\hat{\mu}_k - \hat{\mu}_{\text{tot}})^\top.$$

Set $\hat{Q}^{[m]} := \hat{\Sigma}_w^{-1} \hat{\Sigma}_b$.

end

Set

$$\hat{m} \in \operatorname{argmin}_{m \in [M]} \operatorname{median}(\|\hat{Q}^{[m]} - \hat{Q}^{[m']}\|_{\text{op}} : m' \in [M] \setminus \{m\}). \quad (6)$$

Output: $\hat{Q} := \hat{Q}^{[\hat{m}]}$.

3.2 Theory for base learning using labeled data

In this subsection, we demonstrate how the high-level result in Theorem 2 can be used to derive performance guarantees for a high-dimensional classification algorithm that uses the Sharp-SSL procedure in Algorithm 1 in conjunction with the low-dimensional base method described in Algorithm 2 for estimating the projected whitened between-class covariance matrix. The following theorem provides uniform control of the output of Algorithm 2 for all axis-aligned d -dimensional projected datasets.

Theorem 3. Fix $\varepsilon \in (0, 1]$, $K \in \{2, 3, \dots\}$, $p, d \in \mathbb{N}$ with $p \geq d$ and $n \in \mathbb{N}$ with $n \geq Kd + 1$. Suppose that $(X_1, Y_1), \dots, (X_n, Y_n)$ are independent and identically distributed pairs, with $\mathbb{P}(Y_1 = k) = \pi_k$ and $X_1 | Y_1 = k \sim \mathcal{N}_p(\nu_k^*, \Sigma_w)$ for $k \in [K]$, and let $\psi((PX_i, Y_i)_{i \in [n]})$ be the output of Algorithm 2 with input $(X_i, Y_i)_{i \in [n]}$, for $P \in \mathcal{P}_d$. Suppose that $\|\nu_k^* - \nu_\ell^*\| \leq R_1$ for all $k, \ell \in [K]$ and some $R_1 > 0$, and that Σ_w is diagonal and well-conditioned in the sense that $\max\{\|\Sigma_w\|_{\text{op}}, \|\Sigma_w^{-1}\|_{\text{op}}\} \leq R_2$ for some $R_2 \geq 1$. If

$$\frac{16R_2^2 K}{n} \leq 1 \quad \text{and} \quad \frac{32R_2^2}{n^{1/2}} \log^{1/2} \left(\frac{8 \cdot 9^d \binom{p}{d}}{\varepsilon} \right) \leq 1, \quad (9)$$

then with probability at least $1 - \varepsilon$, we have

$$\begin{aligned} \max_{P \in \mathcal{P}_d} \left\| \psi((PX_i, Y_i)_{i \in [n]}) - (P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top \right\|_{\text{op}} \\ \lesssim_{R_1, R_2} \frac{K}{n} + \sqrt{\frac{d \log(ep/d) + \log(1/\varepsilon)}{n}}. \end{aligned}$$

The sample size condition (9) can be restated as $n \gtrsim_{R_1, R_2} d \log p + \log(1/\varepsilon) + K$, so may be regarded as mild. Regarding K as a constant, Theorem 3 confirms that the uniform control of Algorithm 2 is at the parametric rate, up to a logarithmic factor. The following corollary then follows immediately by combining Theorems 2 and 3.

Corollary 4. Fix $\varepsilon \in (0, 1]$. Suppose that the conditions of Theorem 3 hold, and moreover that $\lambda_{\min}(\Sigma_b) \geq 1/R_3$ for some $R_3 > 0$. Then there exist $C_1, C_2 > 0$, depending only on R_1, R_2 and R_3 , such that if

$$C_1 \left(\frac{K}{n} + \sqrt{\frac{d \log(ep/d) + \log(1/\varepsilon)}{n}} \right) \leq \frac{1}{K},$$

then the output \hat{S} of Algorithm 1 with input $K, p, d \geq s_0, \ell \geq s_0, (X_1, Y_1), \dots, (X_n, Y_n), A, B$, and base procedure ψ from Algorithm 2 satisfies

$$\mathbb{P}(S_0 \subseteq \hat{S}) \geq 1 - \varepsilon - p \exp\left(-\frac{A}{C_2 p^2}\right).$$

Thus, under the conditions of Corollary 4, the Sharp-SSL algorithm can, with high probability, select the signal variables in the top s_0 output variables, provided that the number A of groups of random projections is large by comparison with p^2 . In other words, the algorithm reduces the problem to a low-dimensional one, for which standard learning techniques can be applied. The guarantees for these methods (e.g. Anderson, 2003, Theorem 6.6.1) can then be combined on the high-probability event of Corollary 4 to establish theoretical results for the full procedure.

3.3 Theory for semi-supervised based learning

When the proportion of labeled data is low, Algorithm 2 may be inaccurate when used as the base procedure in Algorithm 1. The aim of this subsection, therefore, is to study the base procedure of Algorithm 3, which is able to leverage both the labeled and unlabeled data via an EM algorithm to estimate the whitened between-class covariance matrix for each projected data set. Our analysis builds on several recent breakthroughs in our understanding of the EM algorithm. This line of work includes Balakrishnan, Wainwright and Yu (2017), Daskalakis, Tzamos and Zampetakis (2017), Yan, Yin and Sarker (2017), Dwivedi et al. (2020a), Dwivedi et al. (2020b), Davis, Diaz and Wang (2021), Ho et al. (2020), Ndaoud (2022), Wu and Zhou (2022) and Doss et al. (2023), all of which focus on the unsupervised case.

For simplicity, we will focus on the setting where independent and identically distributed $(X_1, Y_1^*), \dots, (X_n, Y_n^*)$ are generated from a mixture of two Gaussians with opposite means and identity covariance matrix:

$$Y_i^* \sim \text{Unif}(\{1, 2\}), X_i | Y_i^* \sim \mathcal{N}_p((-1)^{Y_i^*} \nu^*, I_p), \text{ and } Y_i = Y_i^* \mathbb{1}_{\{i \leq n_L\}} \text{ for all } i \in [n]. \quad (10)$$

We assume that we observe $(X_1, Y_1), \dots, (X_{n_L}, Y_{n_L}), X_{n_L+1}, \dots, X_n$ for some $n_L \in \{0, \dots, n\}$. In other words, we are given n_L labeled observations and $n_U := n - n_L$ unlabeled ones. Thus, $n_L = 0$ corresponds to the fully unsupervised case, i.e., clustering, while $n_L = n$ corresponds to the supervised case, i.e., classification. We define $Y_i = Y_i^*$ for $i \in [n_L]$, and $Y_i = 0$ for $i \in \{n_L + 1, \dots, n\}$.

We first study the performance of the EM procedure after the covariates have been projected into a lower-dimensional space. In other words, for some fixed $P \in \mathcal{P}_d$, define $Z_i := PX_i$ for $i \in [n]$ and $\mu^* := P\nu^* \in \mathbb{R}^d$, so that $Z_i | Y_i^* \sim \mathcal{N}_d((-1)^{Y_i^*} \mu^*, I_d)$. In this setting, we have a single unknown parameter μ^* to estimate, and this can be achieved by applying Algorithm 3 to $(Z_i, Y_i)_{i \in [n]}$ with $K = 2$ and the constraint set

$$\mathcal{C} := \{(-\mu, \mu, I_d) : \mu \in \mathbb{R}^d\}. \quad (11)$$

After initializing the EM algorithm at some fixed $(-\hat{\mu}^{(0)}, \hat{\mu}^{(0)}, I_d) \in \mathcal{C}$, for $t \in \mathbb{N}$, the t th iterate of the EM iteration described in (4) and (5) is $(-\hat{\mu}^{(t)}, \hat{\mu}^{(t)}, I_d)$, where

$$\hat{\mu}^{(t)} := \frac{1}{n} \left\{ \sum_{i: Y_i \neq 0} (-1)^{Y_i} Z_i + \sum_{i: Y_i = 0} Z_i \tanh \langle Z_i, \hat{\mu}^{(t-1)} \rangle \right\}; \quad (12)$$

see Lemma 15. Since we allow $n_L = 0$, where μ is only identifiable up to sign, and since the between-class sample covariance matrix $\hat{\Sigma}_b$ computed in Algorithm 3 is equal to $\hat{\Sigma}_b = \hat{\mu}_1 \hat{\mu}_1^\top - \hat{\mu}_{\text{tot}} \hat{\mu}_{\text{tot}}^\top$, which is invariant to flipping the signs of $\hat{\mu}_1$ and $\hat{\mu}_2$ simultaneously, it is natural to consider the loss function $L : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, \infty)$ given by

$$L(\mu, \mu') := \|\mu - \mu'\| \wedge \|\mu + \mu'\|.$$

Proposition 5 below provides a theoretical guarantee for this semi-supervised EM algorithm. For notational simplicity, we define $\gamma := n_L/n$, $\omega_0 := \sqrt{\{d \log n + \log(1/\delta)\}/n_U}$ and $\zeta_0 := \min\{\omega_0 \gamma^{-1/2}, \omega_0^{1/2}\}$ throughout this section. Thus, treating d as a constant and ignoring polylogarithmic terms, ω_0 is of order $n_U^{-1/2}$ and ζ_0 is of order $\min\{n_L^{-1/2}, n_U^{-1/4}\}$ when $\gamma < 1/2$. We remark that $n_L^{-1/2}$ is the critical ℓ_2 -testing radius for distinguishing the means of two labeled Gaussian distributions with identity

covariance using n_L observations. On the other hand, as we show in Lemma 16, no test of the null hypothesis $H_0 : \mathcal{N}_d(0, I_d)$ against the two-component mixture alternative $H_1 : \frac{1}{2}\mathcal{N}_d(\mu^*, I_d) + \frac{1}{2}\mathcal{N}_d(-\mu^*, I_d)$ based on n_U observations can have large power unless the signal strength $\|\mu^*\|$ is at least of order $n_U^{-1/4}$.

Proposition 5. *Fix $\delta \in (2e^{-n}, 1]$ and $r \geq 1$, and suppose that $\|\mu^*\| \leq r$ and $\gamma < 1/2$. There exists $c > 0$, depending only on r , such that if $\omega_0 \leq c$ and $n \geq 3$, then the following statements hold:*

(i) *For any $\hat{\mu}^{(0)} \in \mathbb{R}^d$ with $\|\hat{\mu}^{(0)}\| \leq r + 3$, we have with probability at least $1 - 2\delta$ that*

$$\limsup_{t \rightarrow \infty} L(\hat{\mu}^{(t)}, \mu^*) \lesssim_r \zeta_0 \vee \|\mu^*\|.$$

(ii) *There exists $C > 0$, depending only on r , such that if $\|\mu^*\| \geq C\zeta_0\sqrt{d \log n}$ and $\hat{\mu}^{(0)} = (\zeta_0 \vee r\omega_0)\eta_0$ with $\eta_0 \sim \text{Unif}(\mathbb{S}^{d-1})$, then with probability at least $1 - 2\delta - \sqrt{2}/(\pi \log n_U)$, we have*

$$\limsup_{t \rightarrow \infty} L(\hat{\mu}^{(t)}, \mu^*) \lesssim_r \frac{\omega_0}{\|\mu^*\|} \wedge \frac{\omega_0}{\gamma^{1/2}}.$$

In order to interpret Proposition 5(i), consider the regime where $\|\mu^*\| \leq \zeta_0$. In this case, as discussed above, the two mixture components are essentially indistinguishable, and the bound reveals that the EM algorithm performs no worse than the trivial zero estimator, up to constant factors. On the other hand, part (ii) studies the more interesting regime where the two mixture components are distinguishable, and we establish a faster convergence rate for the EM algorithm in this strong signal regime.

The following theorem combines the two convergence regimes in Proposition 5 to derive a convergence guarantee for the estimated whitened between-class covariance matrix output by Algorithm 3. To state the result, recall the definition of \mathcal{C} from (11). For any $\zeta > 0$, we write $U(\zeta)$ for the pushforward measure on \mathcal{C} induced by $\text{Unif}(\zeta\mathbb{S}^{d-1})$ under the map $\mu \mapsto (-\mu, \mu, I_d)$.

Theorem 6. *Fix $\delta \in (2e^{-n}, 1]$, and $r \geq 1$ and suppose that $\|\mu^*\| \leq r$ and $\gamma < 1/2$. There exists $c > 0$, depending only on r , such that if $\omega_0 \leq \min\{c, (d \log n)^{-3}\}$ and $n \geq 108$, then the sequence of outputs $(\hat{Q}^{(T)})_{T \in \mathbb{N}}$ of Algorithm 3 with inputs $(Z_1, Y_1), \dots, (Z_n, Y_n)$, \mathcal{C} , $\pi_{\mathcal{C}} = U(\zeta_0 \vee r\omega_0)$, $M \in \mathbb{N}$ and $T \in \mathbb{N}$ satisfies with probability at least $1 - 3\delta - e^{-M/50}$ that*

$$\limsup_{T \rightarrow \infty} \|\hat{Q}^{(T)} - \mu^* \mu^{*\top}\|_{\text{op}} \lesssim_r \frac{\omega_0}{\|\mu^*\|} \wedge \zeta_0.$$

Finally in this section, we study the implications of Theorem 6 for the recovery of the signal coordinates in the semi-supervised learning setting. We write $\psi^{(M, T)}$ for the base procedure that takes $(z_i, y_i)_{i \in [n]} \in \mathbb{R}^d \times ([K] \cup \{0\})$ as input and returns the output of Algorithm 3 when run with these inputs together with \mathcal{C} , $\pi_{\mathcal{C}}$, M and T . Let S_0 denote the set of coordinates where $\nu^* \in \mathbb{R}^p$ is non-zero, and let $s_0 := |S_0|$.

Corollary 7. *Fix $\varepsilon \in (8e^{-n/2}, 1]$, $r \geq 1$, and suppose that $\|\mu^*\| \leq r$, $M \geq 50 \log(4/\varepsilon) + 50d \log p$ and $\gamma < 1/2$. Let $\nu_{\max}^* := \|\nu^*\|_{\infty}$ and let ν_{\min}^* denote the minimum absolute value of a non-zero component of ν^* . There exist $C_1, C_2 > 0$, depending only on r , such that if $n \geq C_1(d \log p)^6 \{d \log p + \log(1/\varepsilon)\}$, and*

$$C_2 \min \left[\left\{ \frac{d \log(p \vee n) + \log(1/\varepsilon)}{n} \right\}^{1/4}, \sqrt{\frac{d \log(p \vee n) + \log(1/\varepsilon)}{n_L}} \right] \leq \frac{(\nu_{\min}^*)^2}{4},$$

then the sequence of outputs $(\hat{S}^{(T)})_{T \geq 1}$ of Algorithm 1 with inputs $K = 2$, p , $d \geq s_0$, $\ell \geq s_0$, $(X_i, Y_i)_{i \in [n]}$, A , B and base procedure $\psi^{(M, T)}$ satisfies

$$\liminf_{T \rightarrow \infty} \mathbb{P}(S_0 \subseteq \hat{S}^{(T)}) \geq 1 - \varepsilon - pe^{-A(\nu_{\min}^*)^4 / (50p^2(\nu_{\max}^*)^4)}.$$

Corollary 7 reveals in particular that, treating ν_{\max}^* and ν_{\min}^* as constants and under the stated sample size conditions, we again recover all of the signal coordinates in the top s_0 output entries, provided that A is large by comparison with p^2 . Thus, in this sense, we can achieve a similar guarantee to that provided by Corollary 4, though the number of groups of projections required for a high probability guarantee in Corollary 7 may be significantly larger in settings where the ratio $\nu_{\max}^* / \nu_{\min}^*$ is large.

4 Numerical studies

Throughout this section, unless otherwise stated, data $(X_i, Y_i, Y_i^*)_{i \in [n]}$ are sampled from an equal-probability normal mixture as follows: $\mathbb{P}(Y_i^* = k) = 1/K$ for $k \in [K]$, $\mathbb{P}(Y_i = Y_i^*) = 1 - \mathbb{P}(Y_i = 0) = \gamma$ and $X_i | Y_i^* \sim N_p(\mu_{Y_i^*}, \Sigma_w)$. The cluster means $(\mu_k)_{k \in [K]}$ are chosen to be s_0 -sparse and we define the signal-to-noise ratio of the problem to be³

$$\text{SNR} := \frac{\min_{k, k' \in [K], k \neq k'} \|\mu_k - \mu_{k'}\|}{\sqrt{\text{tr}(\Sigma_w)/p}}. \quad (13)$$

In our numerical studies, we slightly modify Algorithm 3 so that instead of randomly initializing the cluster means and the covariance matrix, we use the output of hierarchical clustering to initialize the EM algorithm as implemented in the `mclust` R package (Fraley and Raftery, 1998). This allow us to run Algorithm 3 with $M = 1$.

4.1 Choice of tuning parameters

The purpose of this subsection is to investigate the effect of the various input parameters A , B , d and ℓ in Algorithm 1, and to recommend sensible default choices. In Figure 1, we plot the misclustering rate with Algorithm 3 as a base procedure in our Gaussian semi-supervised learning setting as each of these parameters varies, for four different SNR levels. After applying Algorithm 1, we obtain our final estimated cluster labels by using Algorithm 3 again on the data projected onto the selected coordinates with a single hierarchical clustering initialization. We then output the predicted labels, computed as $\hat{y}_i := \text{sargmax}_{k \in [K]} L_{i,k}$ for $i \in [n]$, instead of \hat{Q} .

The panels in Figure 1 reveal that the misclustering rate is quite robust to the choices of A , B and d , and that it is less serious (and may even help) to choose ℓ larger—rather than smaller—than s_0 . In particular, it seems that $A = 150$ suffices for almost optimal performance (though there appears to be some penalty for choosing it to be as small as 50), and $B = 75$ appears adequate. There is no clear trend on performance with the choice of d , so for simplicity we took $d = s_0$ in our remaining simulations below. Finally, the misclustering rate appears to decrease as ℓ increases, with an elbow in the curve visible at the highest value of the SNR when ℓ is set to the true sparsity level s_0 . Of course, if ℓ is chosen to be very large, then we will include many noise variables, and

³In some of our simulations, Σ_w was generated randomly for convenience. In such settings, we replaced $\text{tr}(\Sigma_w)/p$ in the denominator of (13) with $\mathbb{E}\{\text{tr}(\Sigma_w)\}/p$.

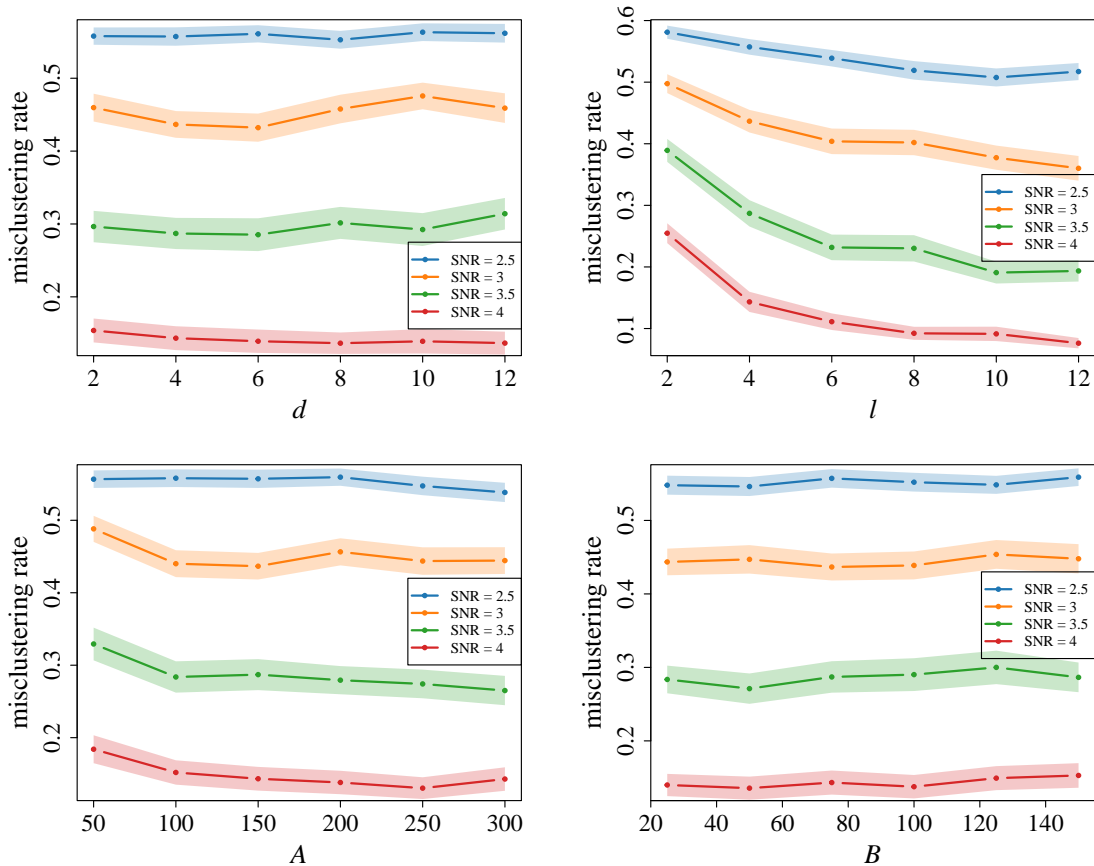


Figure 1: Average misclustering rate over 200 repetitions in our anisotropic Gaussian semi-supervised learning problem with $n = 250$, $p = 600$, $s = 4$, $K = 3$, $\gamma = 0.05$, $\text{SNR} \in \{2.5, 3, 3.5, 4\}$ and $\Sigma_w = V\Lambda V^\top$, where $\Lambda \in \mathbb{R}^{p \times p}$ is diagonal with independent $\text{Unif}[0, 2]$ diagonal entries and V is independent of Λ , and generated according to the Haar measure on $\mathbb{O}^{p \times p}$. For each of the four panels, we fix three of $d = 4$, $l = 4$, $A = 150$, $B = 75$, and vary the remaining one. The shaded regions represent interpolated 95% confidence intervals at each of the points.

the misclustering rate will eventually deteriorate. Nevertheless, the bottom-right panel of Figure 1 indicates that the gain in increasing the probability of including all signal variables may outweigh the penalty of also including more noise variables—as expected, this effect is larger when the SNR is larger. For simplicity we choose $l = s_0$ in our remaining simulations, though we recommend practitioners err on the side of choosing larger l .

4.2 Comparison with existing methods

In this subsection, we compare the empirical performance of the Sharp-SSL algorithm in high-dimensional clustering tasks with several existing approaches. We apply the Sharp-SSL algorithm using the EM algorithm of Algorithm 3 as a base procedure, with input parameters $A = 150$, $B = 75$, $d = l = s_0$ as discussed in Section 4.1, and our final estimated cluster labels are then obtained as described there.

We compare the Sharp-SSL algorithm with five alternative high-dimensional clustering methods: spectral clustering (e.g. von Luxburg, 2007), the ℓ_1 -penalized approach of

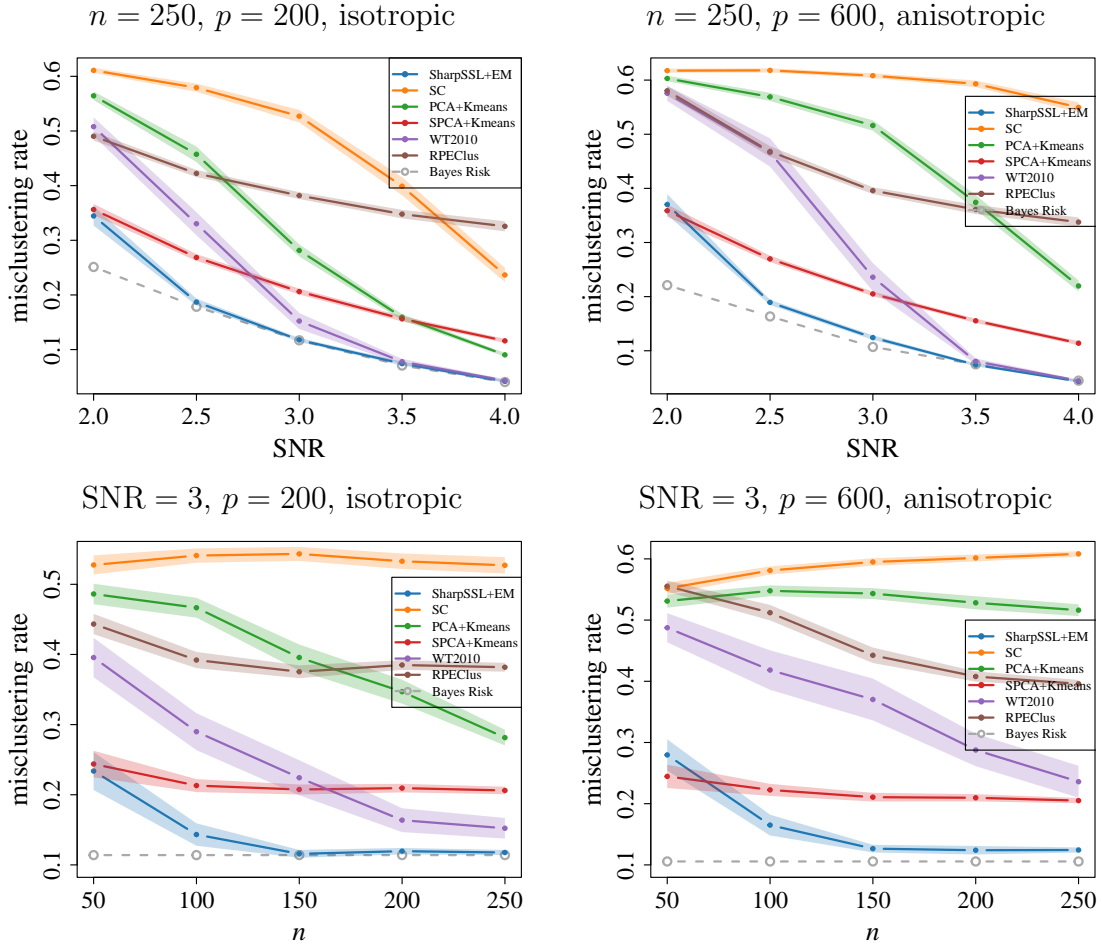


Figure 2: Average misclustering rate over 100 repetitions using **Sharp-SSL** followed by the EM algorithm, as well as using the other methods from Section 4.2. Data are generated from the normal mixture distribution described at the beginning of Section 4 with $K = 3$ and $p = 200$ (left) as well as $p = 600$ (right). The three cluster means are given by $\mu_1 = a(1, 1, 0, \mathbf{0}_{p-3})$, $\mu_2 = a(-1, 0, 1, \mathbf{0}_{p-3})$ and $\mu_3 = a(0, -1, -1, \mathbf{0}_{p-3})$, where the scale a is chosen such that their pairwise distances are all equal to SNR. For isotropic settings (left), $\Sigma_w = I_p$; for anisotropic settings (right), $\Sigma_w = V\Lambda V^\top$, where $\Lambda \in \mathbb{R}^{p \times p}$ is diagonal with independent $\text{Unif}[0, 2]$ diagonal entries and V is independent of Λ , and sampled from the Haar measure on $\mathbb{O}^{p \times p}$. The Bayes risk is shown as the gray dashed line. In the top panels, $n = 250$ and the SNR varies; in the bottom panels, SNR = 3 and n varies. The shaded regions represent interpolated 95% confidence intervals at each of the points.

Witten and Tibshirani (2010) and the RPECLUS algorithm of Anderlucci, Fortunato and Montanari (2022) as well as a pair of methods that, like Sharp-SSL, apply dimension reduction prior to a low-dimensional clustering algorithm.

In more detail, the spectral clustering approach first constructs a J -nearest neighbour graph adjacency matrix $A = (A_{i,i'})_{i,i' \in [n]} \in \{0, 1\}^{n \times n}$, where $A_{i,i'} := 1$ if either X_i is one of the $J = 10$ nearest neighbours of $X_{i'}$ in Euclidean distance or vice versa, and $A_{i,i'} := 0$ otherwise. It then computes an $n \times K$ matrix of eigenvectors associated with the K smallest nonzero eigenvalues of the Laplacian matrix $L := D - A$, where $D \in \mathbb{R}^{n \times n}$ is a diagonal matrix with diagonal entries $D_{i,i} := \sum_{i' \in [n]} A_{i,i'}$. The final step is to apply the K -means clustering algorithm (Lloyd, 1982), as implemented in the `kmeans` base R function with 100 random initializations, to the rows of L with the oracle choice of K .

The Witten and Tibshirani (2010) method, which is implemented in the `sparcl` R package, determines the estimated cluster memberships by maximizing a coordinatewise-weighted between-cluster sum of squares criterion, subject to an ℓ_1 constraint on the weights. A permutation approach is used to select the ℓ_1 tuning parameter.

In the RPECLUS algorithm of Anderlucci, Fortunato and Montanari (2022), we generate B random orthogonal projections and incorporate the d -dimensional projected data as covariates for a linear regression with the orthogonal complement of the projected data as the response. We then use the Bayesian Information Criteria (BIC) from both an application of the EM algorithm to the projected data and the aforementioned regression to identify good projections, and aggregate using the consensus clustering technique of Dimitriadou, Weingessel and Hornik (2002) over the best B^* projections chosen according to the sum of the BIC scores. Following the recommendation of Anderlucci, Fortunato and Montanari (2022), we took $B = 1000$ and $B^* = 100$ as well as $d = s_0$. It turned out that this approach had a misclustering rate almost identical to that of a random guess, likely because it did not leverage the sparsity of the signal. We therefore modified this method by generating random axis-aligned projections instead of orthogonal ones, and report this version in our comparison.

The first of the two-stage approaches applies principal component analysis (PCA) to project the data into the oracle choice of $K - 1$ dimensions (the dimension of the space spanned by the K cluster means); the second uses sparse principal component analysis (SPCA), as implemented in the SPCA_{VRP} algorithm (Gataric, Wang and Samworth, 2020) with inputs $A = 600$, $B = 200$, and the oracle choices $d = \ell = s_0$, to project into s_0 dimensions. Thereafter, both algorithms apply K -means to the projected data as above. We also explored the option of replacing the K -means steps in these latter algorithms with the EM algorithm, but observed very little difference, so do not report these results here.

Given true labels $y_1, \dots, y_n \in [K]$ and estimated labels $\hat{y}_1, \dots, \hat{y}_n \in [K]$ from a clustering algorithm, we measure the performance of the algorithm via its *misclustering rate*, defined as⁴

$$L(\{y_1, \dots, y_n\}, \{\hat{y}_1, \dots, \hat{y}_n\}) := \min_{\sigma \in S_K} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\sigma(\hat{y}_i) \neq y_i\}},$$

where S_K is the group of all permutations of $[K]$. In particular, Figure 2 presents the average misclustering rates over 100 Monte Carlo repetitions of the different high-dimensional clustering algorithms described above. Across two different dimensions $p \in \{200, 600\}$,

⁴Here, the minimum over permutations is taken because it is only the cluster groupings, and not the labels themselves, that are important.

isotropic and anisotropic settings, and for different values of n and SNR, we see a consistent picture of the **Sharp-SSL** algorithm combined with EM producing the lowest misclustering rates, often by a large margin. Indeed, for all but the smallest sample sizes or values of SNR, the **Sharp-SSL+EM** algorithm nearly attains the Bayes risk in all of the problems considered here.

4.3 Effect of observed fraction on misclustering rate

One of the key attractions of our procedure is that it offers a unified framework to perform classification or clustering with an arbitrary fraction of labeled observations. In this subsection, we explore the performance of the algorithm as we vary the proportion of observed labels.

Recall that we have two different options for the way in which we implement the **Sharp-SSL** algorithm to estimate the set of signal coordinates: we can either use only the labeled data, as in the supervised learning approach of Algorithm 2, or we can try to leverage in addition the unlabeled data via the semi-supervised EM approach of Algorithm 3. In the extreme case of this latter version, we have no labeled data, so the algorithm is unsupervised. In Figure 3 we compare the performance of these three methods in both high- and low-dimensional versions of the normal mixture distribution data generation mechanism described at the beginning of Section 4 as the proportion γ of observed labels varies.

More precisely, for the semi-supervised and unsupervised algorithms, we adopt the same implementation of **Sharp-SSL** as described at the beginning of Section 4.2. The supervised algorithm is very similar, but applies Algorithm 2 in place of Algorithm 3 to select coordinates, and obtains final predicted labels by applying LDA again on the projected labeled data. In cases where the proportion of labeled data was so small that the convex hull of the projected labeled data was less than full-dimensional for every class, we forced Algorithm 2 to return a zero matrix (this only happened when γ was very small).

The top panels of Figure 3 present the results in high-dimensional settings with $p \in \{200, 600\}$. Since the unsupervised approach has no access to the labels, it has constant misclustering rate. The performance of the semi-supervised approach is always at least as good as that of the unsupervised algorithm, and improves as γ increases. In other words, it effectively leverages the additional information provided by the class labels. When γ is very small, the supervised algorithm—which ignores the unlabeled data—is inaccurate, as it has very little data to work with. On the other hand, its performance also improves as γ increases, and once around 5% of our data are labeled, it outperforms the unsupervised algorithm. Further, it essentially matches the semi-supervised approach when about a third of the data are labeled. We truncate the plot at $\gamma = 1/2$ to ensure that we have enough test data on which to compute the misclustering rate.

In the bottom panels of Figure 3, we explore the performance of the three algorithms above in two low-dimensional settings with different values of SNR, in order to provide further insight into the phenomena described in the previous paragraph. Here, we take $K = 2$ and report the average Frobenius norm loss

$$\mathcal{L}((\hat{\mu}_1, \hat{\mu}_2), (\mu_1, \mu_2)) := \min \left\{ \|(\hat{\mu}_1, \hat{\mu}_2) - (\mu_1, \mu_2)\|_F, \|(\hat{\mu}_2, \hat{\mu}_1) - (\mu_1, \mu_2)\|_F \right\}$$

of the estimated means, over 100 repetitions. If there are insufficient labeled data to run Algorithm 2, then we output $\hat{\mu}_1 = \hat{\mu}_2 = \mathbf{0}_p$. We see that, already in these low-

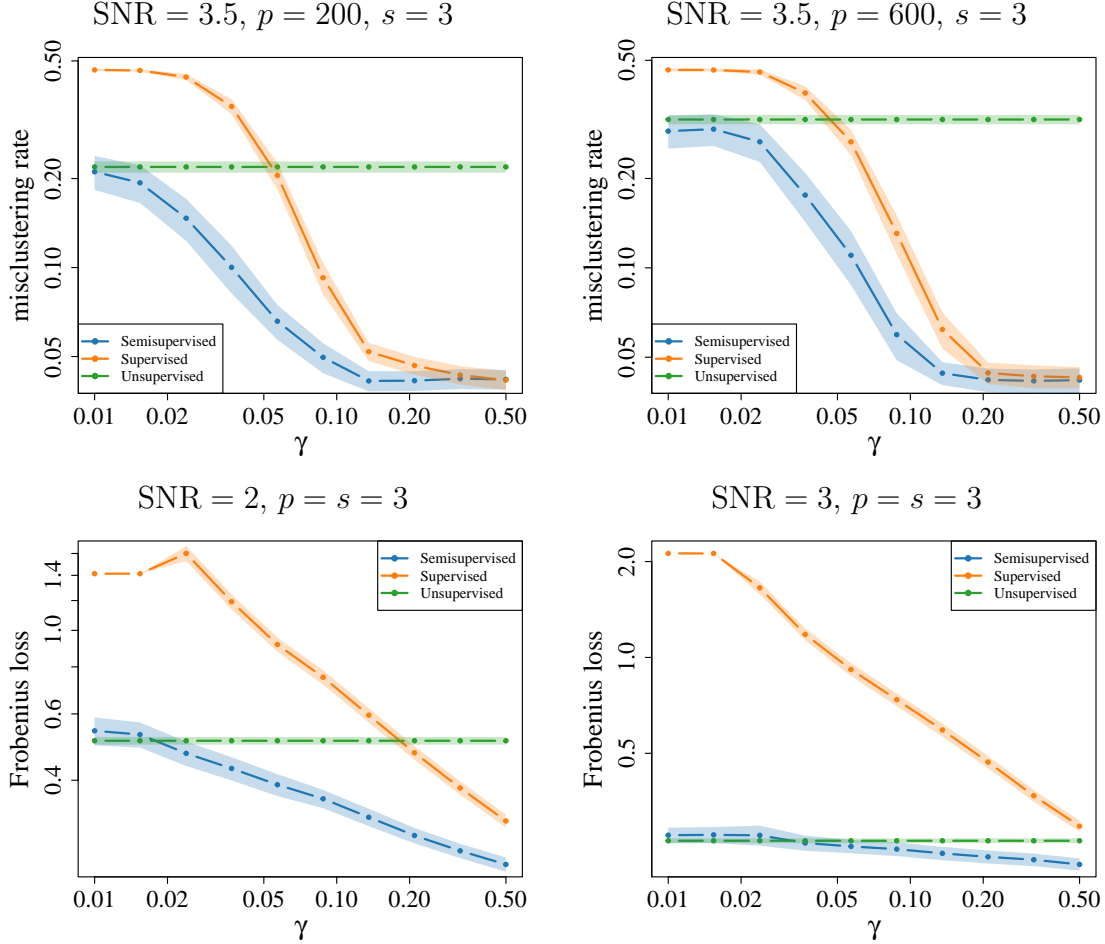


Figure 3: Effect of label fraction on performance of supervised, semi-supervised and unsupervised Sharp-SSL learning methods. Data are generated from the normal mixture distribution described at the beginning of Section 4 with $K = 2$ and $\Sigma_w = I_p$, $\mu_1 = -\mu_2 = a(\mathbf{1}_s, \mathbf{0}_{p-s})^\top \in \mathbb{R}^p$, where a is chosen such that $\|\mu_1 - \mu_2\| = \text{SNR}$. Bottom: average Frobenius loss of estimating the $(\mu_1, \mu_2) \in \mathbb{R}^{p \times 2}$ over 100 repetitions via the semi-supervised approach (Algorithm 3), supervised approach (Algorithm 2) and unsupervised approach (Algorithm 3 without using the labels). Top: average misclustering rate over 100 repetitions from applying the above three methods as base algorithms in Algorithm 1. The shaded regions represent interpolated 95% confidence intervals at each of the points.

dimensional problems, a similar picture emerges: if the proportion of labeled data is small, then the unsupervised algorithm outperforms the supervised one, but this situation may be reversed when γ is larger. The semi-supervised algorithm is able to leverage both the unlabeled and labeled data to obtain the best of both worlds. These empirical observations agree with our theory from Section 3, in particular in the way in which Theorem 6 bounds the accuracy of mean estimation for the semi-supervised algorithm by a minimum of a term that does not depend on γ and one that decreases as γ increases. It appears that the switch in the minimum occurs around $\gamma = 0.02$ in these examples.

4.4 Empirical data analysis

In this subsection we apply the **Sharp-SSL** algorithm, as well as several competing methods, to the gene expression data set from Alon et al. (1999), which contains observations on 62 patients. A preprocessed version of the data can be downloaded from the R package ‘datamicroarray’ (Ramey, 2016), with a total of 2000 features (genes) measured on 40 patients with colon tumors and 22 without tumors. We first exclude 9 genes to remove perfect collinearity and then standardize each of the remaining $p = 1991$ columns of the dataset to have unit variance.

We apply the **Sharp-SSL** algorithm using EM (Algorithm 3) as the base procedure, with input parameters $A = 150$, $B = 75$, $d = \ell = 5$. In addition to our approach (**Sharp-SSL+EM**), we also compare the performance of the spectral clustering (SC) method, the Witten and Tibshirani (2010) method (WT2010), as well as four two-stage methods (PCA+Kmeans, PCA+EM, SPCA+Kmeans, SPCA+EM), where we first reduce dimension of the data to a 5-dimensional subspace using either PCA or SPCA and then apply either the EM algorithm or K -means clustering on the low-dimensional data. For SPCA, we use the SPCAvRP algorithm (Gataric, Wang and Samworth, 2020) with inputs $A = 600$, $B = 200$ and $d = \ell = 5$. The true labels are hidden to all algorithms and are only used to evaluate the final misclustering rate.

Over 100 Monte Carlo repetitions of the randomized algorithms, the **Sharp-SSL+EM** method had an average misclustering rate of 28.8%, whereas all other competitors had a misclustering rate above 40%, as can be seen from the right-hand data points in Figure 4. To investigate this performance further, we applied each method to a subset of the features. These were constructed from the top $\ell = 5$ genes identified through **Sharp-SSL**, together with $m = 0, 10, 50, 200$ and 600 randomly chosen genes from the remaining 1986. The results are presented as the other data points in Figure 4. We see that the improved performance of the **Sharp-SSL+EM** relative to the other methods persists, even when only a small number of potentially non-discriminative covariates are present. When $m = 0$, **Sharp-SSL+EM** has a slight disadvantage as other algorithms benefit from the ensemble effect of combining two different learning methods; nevertheless it remains competitive. This reinforces the point that the primary contribution of the **Sharp-SSL** algorithm is to identify signal coordinates that are helpful for semi-supervised learning, and once this task has been accomplished, a variety of low-dimensional procedures are available to the practitioner.

5 Proof of the main results

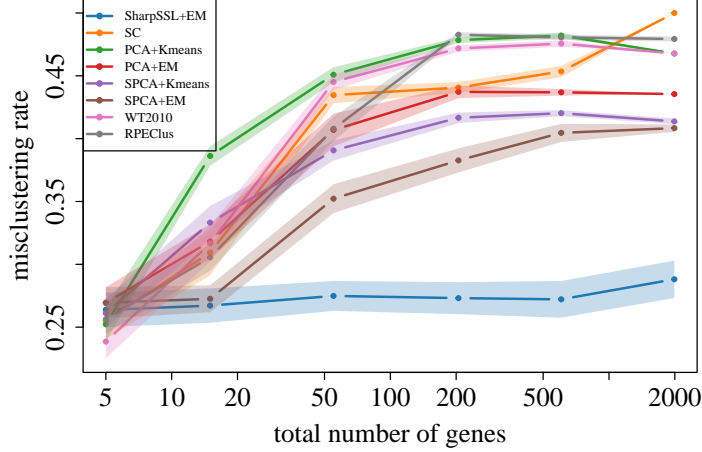


Figure 4: Average misclustering rate (over 100 repetitions for randomized algorithms) for the colon tumor data, using **Sharp-SSL** followed by the EM algorithm, as well as the other methods described in Section 4.4. The right-hand data points plot the average misclustering rate on the full data set. The other points were obtained by applying each method to a subset of genes formed from the top five genes identified by **Sharp-SSL** together with randomly sampled genes. The shaded regions represent interpolated 95% confidence intervals at each of the points.

5.1 Proof of Proposition 1

The assumption that the convex hull of ν_1, \dots, ν_K is $(K - 1)$ -dimensional implies that Σ_b is of rank $K - 1$. Define $A := \Sigma_w^{-1/2} \Sigma_b \Sigma_w^{-1/2} \in \mathbb{R}^{p \times p}$, which has rank $K - 1$. Given $V \in \mathbb{O}^{p \times d}$, we can find $Q \in \mathbb{O}^{p \times d}$ with the same column span as that of $\Sigma_w^{1/2} V$ and let $R := Q^\top \Sigma_w^{1/2} V \in \mathbb{R}^{d \times d}$, so that R is invertible, and $\Sigma_w^{1/2} V = QR$. We observe that

$$\text{tr}\{(V^\top \Sigma_w V)^{-1} (V^\top \Sigma_b V)\} = \text{tr}\{(R^\top R)^{-1} (R^\top Q^\top A Q R)\} = \text{tr}(Q^\top A Q).$$

Thus, $J(V; \Sigma_b, \Sigma_w)$ depends on V only through the column space of $\Sigma_w^{1/2} V$. Moreover, $\text{tr}(Q^\top A Q)$ is maximized when Q , or equivalently $\Sigma_w^{1/2} V$, spans a d -dimensional space that contains the $(K - 1)$ -dimensional eigenspace corresponding to the non-zero eigenvalues of A . Note that if for some $v \in \mathbb{R}^p \setminus \{0\}$ and $\lambda \geq 0$, we have $Av = \lambda v$, then $\Sigma_w^{-1} \Sigma_b \Sigma_w^{-1/2} v = \Sigma_w^{-1/2} Av = \lambda \Sigma_w^{-1/2} v$, so $\Sigma_w^{-1/2} v$ is an eigenvector of $\Sigma_w^{-1} \Sigma_b$ with eigenvalue λ . Hence V maximizes $J(V; \Sigma_b, \Sigma_w)$ over $\mathbb{O}^{p \times d}$ if and only if V spans a d -dimensional space that contains the $(K - 1)$ -dimensional eigenspace corresponding to the $K - 1$ non-zero eigenvalues of $\Sigma_w^{-1} \Sigma_b$. Finally, for any $v \in \mathbb{R}^p \setminus \{0\}$,

$$v^\top \Sigma_w^{-1/2} \Sigma_b \Sigma_w^{-1/2} v = \sum_{k=1}^K \pi_k v^\top \Sigma_w^{-1/2} (\nu_k - \nu) (\nu_k - \nu)^\top \Sigma_w^{-1/2} v \neq 0$$

if and only if $v^\top \Sigma_w^{-1/2} (\nu_k - \nu) \neq 0$ for some $k \in [K]$. Thus, the eigenspace corresponding to the non-zero eigenvalues of A is spanned by $(\Sigma_w^{-1/2} (\nu_k - \nu) : k \in [K])$, and so the eigenspace corresponding to non-zero eigenvalues of $\Sigma_w^{-1} \Sigma_b$ is spanned by $(\Sigma_w^{-1} (\nu_k - \nu) : k \in [K])$.

5.2 Proof of Theorem 2

We write $S^{a,b} := \{j \in [p] : (P^{a,b,\top} P^{a,b})_{j,j} = 1\}$. For any $S = \{j_1, \dots, j_d\}^\top \in \binom{[p]}{d}$, we identify the set S with the sequence $j_{i_1} < \dots < j_{i_d}$ sorted in increasing order; and with a slight abuse of notation, we will use S to refer to either object, which will always be clear depending on the context. We define $P^S \in \mathcal{P}_d$ by $(P^S)_{\ell,j} := \mathbb{1}_{\{j=j_\ell\}}$, so that $P^{S,\top} P^S = \text{diag}((\mathbb{1}_{\{j \in S\}})_{j \in [p]})$. Define $Q^S := (P^S \Sigma_w P^{S,\top})^{-1} P^S \Sigma_b P^{S,\top} \in \mathbb{R}^{d \times d}$ and $\hat{Q}^S := \psi((P^S X_i, Y_i)_{i \in [n]}) \in \mathbb{R}^{d \times d}$. Note that $\hat{Q}^{a,b} = \hat{Q}^{S^{a,b}}$ in this notation, and we will similarly denote $Q^{a,b} := Q^{S^{a,b}}$ for simplicity. Recalling the definition of Ω from (8), in our new notation, and recalling that ψ is permutation-equivariant, we can write

$$\Omega = \left\{ \max_{S \in \binom{[p]}{d}} \|\hat{Q}^S - Q^S\|_{\text{op}} < \frac{\gamma_{\min}}{4(K-1)} \right\},$$

and have $\mathbb{P}(\Omega) \geq 1 - \varepsilon$ by (7). We will work on the event Ω throughout the remainder of the proof, and assume also that $\gamma_{\min} > 0$, because otherwise the conclusion is trivial.

By Weyl's inequality (Weyl, 1912; Stewart and Sun, 1990, Corollary IV.4.9), we have on Ω that

$$|\text{tr}(\hat{Q}^S) - \text{tr}(Q^S)| \leq (K-1) \|\hat{Q}^S - Q^S\|_{\text{op}} \leq \frac{\gamma_{\min}}{4}.$$

On the other hand, we note that $\text{tr}(Q^S) = \sum_{j \in S \cap S_0} (\Sigma_w^{-1} \Sigma_b)_{j,j}$. Therefore, by the triangle inequality, for any $S, S' \in \binom{[p]}{d}$ such that $S \cap S_0$ is a proper subset of $S' \cap S_0$, we have on Ω that

$$\text{tr}(\hat{Q}^S) - \text{tr}(\hat{Q}^{S'}) \leq \frac{\gamma_{\min}}{2} - \sum_{j \in (S' \setminus S) \cap S_0} (\Sigma_w^{-1} \Sigma_b)_{j,j} < 0. \quad (14)$$

Fix $a \in [A]$, and for any $\tilde{j} \in [p]$, define $q_{\tilde{j}} := \mathbb{P}(\tilde{j} \in S^{a,b^*(a)} \mid (X_i, Y_i)_{i \in [n]})$. Now fix some $j \in S_0$ and $j' \in [p] \setminus S_0$. We claim that

$$(q_j - q_{j'}) \mathbb{1}_\Omega \geq 0. \quad (15)$$

To verify this claim, define for $\tilde{j} \in \{j, j'\}$ and $b \in [B]$ the sets

$$\mathcal{S}_{b,\tilde{j}} := \{(S^{a,1}, \dots, S^{a,B}) : b^*(a) = b, \tilde{j} \in S^{a,b}\} \quad \text{and} \quad \mathcal{S}_b := \{(S^{a,1}, \dots, S^{a,B}) : b^*(a) = b\}.$$

Let $f : \binom{[p]}{d} \rightarrow \binom{[p]}{d}$ be a map defined by

$$f(S) := \begin{cases} (S \setminus \{j'\}) \cup \{j\} & \text{if } j \notin S \text{ and } j' \in S \\ S & \text{otherwise.} \end{cases}$$

If $j \notin S^{a,b}$ and $j' \in S^{a,b}$, then $f(S^{a,b}) \cap S_0 = (S^{a,b} \cup \{j\}) \cap S_0 = (S^{a,b} \cap S_0) \cup \{j\}$, so $S^{a,b} \cap S_0$ is a proper subset of $f(S^{a,b}) \cap S_0$; on the other hand, if either $j \in S^{a,b}$ or $j, j' \notin S^{a,b}$, then $f(S^{a,b}) = S^{a,b}$. It follows by (14) that on Ω we have

$$\text{tr}(\hat{Q}^{S^{a,b}}) \leq \text{tr}(\hat{Q}^{f(S^{a,b})}). \quad (16)$$

Now let $F : \mathcal{S}_{b,j'} \rightarrow \mathcal{S}_{b,j}$ be defined as

$$F(S^{a,1}, \dots, S^{a,B}) := (S^{a,1}, \dots, S^{a,b-1}, f(S^{a,b}), S^{a,b+1}, \dots, S^{a,B}).$$

We claim that F is both well-defined and injective on Ω . For the first of these claims, we note that since $j' \in S^{a,b}$, we must have $j \in f(S^{a,b})$. Moreover, if $(S^{a,1}, \dots, S^{a,B}) \in \mathcal{S}_{b,j'}$, then $b^*(a) = b$. But (16) holds on Ω , so $(S^{a,1}, \dots, S^{a,b-1}, f(S^{a,b}), S^{a,b+1}, \dots, S^{a,B}) \in \mathcal{S}_{b,j}$. Hence F is well-defined. For the second claim, suppose that $S_1, S_2 \in \binom{[p]}{d}$ are such that $j' \in S_1 \cap S_2$ and $f(S_1) = f(S_2)$. If $j \in S_1 \cap S_2$, then $S_1 = f(S_1) = f(S_2) = S_2$; if $j \in S_1$ but $j \notin S_2$, then $j \in f(S_2) \setminus f(S_1)$, a contradiction. Similarly, we cannot have $j \notin S_1$ but $j \in S_2$. Finally, if $j \in S_1^c \cap S_2^c$, then $S_1 = (f(S_1) \setminus \{j\}) \cup \{j'\} = (f(S_2) \setminus \{j\}) \cup \{j'\} = S_2$. We deduce that f is injective on $\{S : j' \in S\}$. Since $j' \in S^{a,b}$ for $(S^{a,1}, \dots, S^{a,B}) \in \mathcal{S}_{b,j'}$, this establishes the injectivity of F . In particular, $|\mathcal{S}_{b,j'}| \leq |\mathcal{S}_{b,j}|$. Consequently, on Ω , we have for all $b \in [B]$ that

$$\begin{aligned} \mathbb{P}(j \in S^{a,b^*(a)} \mid (X_i, Y_i)_{i \in [n]}, b^*(a) = b) &= \frac{\mathbb{P}(j \in S^{a,b^*(a)}, b^*(a) = b \mid (X_i, Y_i)_{i \in [n]})}{\mathbb{P}(b^*(a) = b \mid (X_i, Y_i)_{i \in [n]})} = \frac{|\mathcal{S}_{b,j}|}{|\mathcal{S}_b|} \\ &\geq \frac{|\mathcal{S}_{b,j'}|}{|\mathcal{S}_b|} = \frac{\mathbb{P}(j' \in S^{a,b^*(a)}, b^*(a) = b \mid (X_i, Y_i)_{i \in [n]})}{\mathbb{P}(b^*(a) = b \mid (X_i, Y_i)_{i \in [n]})} \\ &= \mathbb{P}(j' \in S^{a,b^*(a)} \mid (X_i, Y_i)_{i \in [n]}, b^*(a) = b), \end{aligned}$$

which implies Claim (15). We remark that one consequence of (15) is that, since $d \geq s_0$, we have on Ω that

$$q_j \geq \frac{\sum_{\tilde{j} \in ([p] \setminus S_0) \cup \{j\}} q_{\tilde{j}}}{p - s_0 + 1} = \frac{d - \sum_{\tilde{j} \in S_0 \setminus \{j\}} q_{\tilde{j}}}{p - s_0 + 1} \geq \frac{d - s_0 + 1}{p - s_0 + 1} \geq \frac{1}{p}. \quad (17)$$

Again fixing $j \in S_0$ and $j' \notin S_0$, we observe on $\Omega \cap \{j \in S^{a,b^*(a)}\}$ that

$$\begin{aligned} \frac{3}{4} \gamma_{\min} &\leq [P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j,j} - \|\hat{Q}^{a,b^*(a)} - Q^{a,b^*(a)}\|_{\text{op}} \leq [P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j,j} \\ &\leq [P^{a,b^*(a), \top} Q^{a,b^*(a)} P^{a,b^*(a)}]_{j,j} + \|\hat{Q}^{a,b^*(a)} - Q^{a,b^*(a)}\|_{\text{op}} \leq \frac{5}{4} \gamma_{\max}, \end{aligned}$$

and similarly on $\Omega \cap \{j' \in S^{a,b^*(a)}\}$ that $|[P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j',j'}| \leq \gamma_{\min}/4$. Recall also that $[P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{\tilde{j}, \tilde{j}} = 0$ for all $\tilde{j} \notin S^{a,b^*(a)}$. Combining the above bounds on the diagonal entries of $\hat{Q}^{a,b^*(a)}$ with (15) and (17), we have on Ω that

$$\begin{aligned} &\mathbb{E}([P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j,j} - [P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j',j'} \mid (X_i, Y_i)_{i \in [n]}) \\ &= \mathbb{E}([P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j,j} \mathbb{1}_{\{j \in S^{a,b^*(a)}\}} \\ &\quad - [P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j',j'} \mathbb{1}_{\{j' \in S^{a,b^*(a)}\}} \mid (X_i, Y_i)_{i \in [n]}) \\ &\geq \frac{q_j \gamma_{\min}}{2} \geq \frac{\gamma_{\min}}{2p}. \end{aligned} \quad (18)$$

Now, let a, j, j' be freely varying again. Since $\hat{w}_j = A^{-1} \sum_{a \in [A]} [P^{a,b^*(a), \top} \hat{Q}^{a,b^*(a)} P^{a,b^*(a)}]_{j,j}$, on Ω we have for any $j \in S_0$ and $j' \notin S_0$ that $\mathbb{E}(\hat{w}_j - \hat{w}_{j'} \mid (X_i, Y_i)_{i \in [n]}) \geq \gamma_{\min}/(2p)$

from (18). Since $\ell \geq s_0$, we have by Hoeffding's inequality that on Ω ,

$$\begin{aligned} \mathbb{P}(S_0 \not\subseteq \hat{S} \mid (X_i, Y_i)_{i \in [n]}) &\leq \mathbb{P}\left(\min_{j \in S_0} \hat{w}_j \leq \max_{j' \notin S_0} \hat{w}_{j'} \mid (X_i, Y_i)_{i \in [n]}\right) \\ &\leq \sum_{j \in S_0} \mathbb{P}\left\{\hat{w}_j - \mathbb{E}(\hat{w}_j \mid (X_i, Y_i)_{i \in [n]}) \leq -\frac{\gamma_{\min}}{4p} \mid (X_i, Y_i)_{i \in [n]}\right\} \\ &\quad + \sum_{j \notin S_0} \mathbb{P}\left\{\hat{w}_j - \mathbb{E}(\hat{w}_j \mid (X_i, Y_i)_{i \in [n]}) \geq \frac{\gamma_{\min}}{4p} \mid (X_i, Y_i)_{i \in [n]}\right\} \\ &\leq p \exp\left\{-\frac{A}{2} \left(\frac{\gamma_{\min}}{4p}\right)^2 \middle/ \left(\frac{5\gamma_{\max}}{4}\right)^2\right\} \leq pe^{-A\gamma_{\min}^2/(50p^2\gamma_{\max}^2)}, \end{aligned}$$

as desired.

5.3 Proof of Theorem 3

The main ingredient of the proof of Theorem 3 is the following proposition, which controls the rate of convergence of the sample between- and within-class covariance matrices to their respective population versions in a classification problem.

Proposition 8 (Rate of convergence for LDA). *Suppose that $(Z_1, Y_1), \dots, (Z_n, Y_n) \in \mathbb{R}^d \times [K]$ are independent and identically distributed data-label pairs, such that $\mathbb{P}(Y_1 = k) = \pi_k$ and $Z_1 \mid Y_1 = k \sim \mathcal{N}_d(\mu_k, \Sigma_w)$ for $k \in [K]$. Write $\mu := \sum_{k \in [K]} \pi_k \mu_k$ and $\Sigma_b := \sum_{k \in [K]} \pi_k (\mu_k - \mu)(\mu_k - \mu)^\top$ and let $\hat{\Sigma}_w$ and $\hat{\Sigma}_b$ be computed as in (3) applied to $(Z_1, Y_1), \dots, (Z_n, Y_n)$. If $\|\mu_k - \mu\| \leq R_1$ for all $k \in [K]$ and $\|\Sigma_w\|_{\text{op}} \leq R_2$ for some $R_1, R_2 > 0$, then for every $\delta \in (0, 1/4]$, we have with probability at least $1 - \delta$ that*

$$\begin{aligned} \|\hat{\Sigma}_b - \Sigma_b\|_{\text{op}} &\leq \frac{12R_2\{K + \log(8 \cdot 9^d/\delta)\}}{n} + (4R_1\sqrt{R_2} + R_1^2 + R_1)\sqrt{\frac{2\log(8 \cdot 9^d/\delta)}{n}}, \\ \|\hat{\Sigma}_w - \Sigma_w\|_{\text{op}} &\leq \frac{4R_2\{K + \log(8 \cdot 9^d/\delta)\}}{n} + 4R_2\sqrt{\frac{\log(8 \cdot 9^d/\delta)}{n}}. \end{aligned}$$

Proof. We first control the rate of convergence of $\hat{\Sigma}_b$. For $n_k := \sum_{i=1}^n \mathbb{1}_{\{Y_i=k\}}$ and $\tilde{\mu} := \sum_{k \in [K]} (n_k/n) \mu_k$, we define $\tilde{\Sigma}_b^{(1)} := \sum_{k=1}^K (n_k/n) (\mu_k - \mu)(\mu_k - \mu)^\top$ and $\tilde{\Sigma}_b^{(2)} := \sum_{k=1}^K (n_k/n) (\mu_k - \tilde{\mu})(\mu_k - \tilde{\mu})^\top$. We have the following decomposition

$$\|\hat{\Sigma}_b - \Sigma_b\|_{\text{op}} \leq \|\tilde{\Sigma}_b^{(1)} - \Sigma_b\|_{\text{op}} + \|\tilde{\Sigma}_b^{(2)} - \tilde{\Sigma}_b^{(1)}\|_{\text{op}} + \|\hat{\Sigma}_b - \tilde{\Sigma}_b^{(2)}\|_{\text{op}}. \quad (19)$$

We will control the three terms on the right-hand side above separately. For the first term, we note that $\|\mu_k - \mu\| \leq R_1$ for all $k \in [K]$. Since n_1, \dots, n_K are functions of Y_1, \dots, Y_n , we have by McDiarmid's inequality (see, e.g. [Boucheron, Lugosi and Massart, 2013](#), Theorem 6.2) that with probability at least $1 - \delta/4$,

$$\|\tilde{\Sigma}_b^{(1)} - \Sigma_b\|_{\text{op}} \leq R_1^2 \sum_{k=1}^K \left| \frac{n_k}{n} - \pi_k \right| \leq R_1^2 \sqrt{\frac{2\log(4/\delta)}{n}}. \quad (20)$$

For the second term, we first apply McDiarmid's inequality again to see that with probability at least $1 - \delta/4$, we have

$$\|\tilde{\mu} - \mu\| \leq \sum_{k=1}^K \left| \frac{n_k}{n} - \pi_k \right| \|\mu_k - \mu\| \leq R_1 \sqrt{\frac{2\log(4/\delta)}{n}}.$$

Thus, we have with probability at least $1 - \delta/4$ that

$$\begin{aligned} \|\tilde{\Sigma}_b^{(2)} - \tilde{\Sigma}_b^{(1)}\|_{\text{op}} &\leq 2 \sum_{k=1}^K \frac{n_k}{n} \|(\mu_k - \mu)(\mu - \tilde{\mu})^\top\|_{\text{op}} + \sum_{k=1}^K \frac{n_k}{n} \|(\mu - \tilde{\mu})(\mu - \tilde{\mu})^\top\|_{\text{op}} \\ &\leq 2R_1 \|\mu - \tilde{\mu}\| + \|\mu - \tilde{\mu}\|^2 \leq 3R_1 \|\mu - \tilde{\mu}\| \leq 3R_1^2 \sqrt{\frac{2 \log(4/\delta)}{n}}. \end{aligned} \quad (21)$$

Finally, for the third term, we write $\hat{\mu}_k := \sum_{i:Y_i=k} Z_i$ and note that $V_k := n_k^{1/2} \hat{\mu}_k$ satisfies $V_k | Y_1, \dots, Y_n \sim \mathcal{N}_d(n_k^{1/2} \mu_k, \Sigma_w)$. Defining $N := (n_1^{1/2}, \dots, n_K^{1/2})^\top$ and $P := NN^\top/n \in \mathbb{R}^{K \times K}$, we may write

$$n\hat{\Sigma}_b = \sum_{k=1}^K n_k \hat{\mu}_k \hat{\mu}_k^\top - n \hat{\mu} \hat{\mu}^\top = V^\top (I_K - P) V,$$

where $V := (V_1, \dots, V_K)^\top$, and where $\hat{\mu} := n^{-1} \sum_{i=1}^n Z_i = \sum_{k=1}^K (n_k/n) \hat{\mu}_k$. By Lemma 17, we deduce that $n\hat{\Sigma}_b$ conditional on Y_1, \dots, Y_n has a d -dimensional non-central Wishart distribution with $K - 1$ degrees of freedom, covariance matrix Σ_w and non-centrality matrix $n\tilde{\Sigma}_b^{(2)}$, which we denote as

$$n\hat{\Sigma}_b | Y_1, \dots, Y_n \sim \mathcal{W}_d(K - 1, \Sigma_w; n\tilde{\Sigma}_b^{(2)});$$

a formal definition is given just before Lemma 17. For any fixed $u \in \mathbb{S}^{d-1}$, we have by Muirhead (2009, Theorem 10.3.6) that

$$u^\top \hat{\Sigma}_b u | Y_1, \dots, Y_n \sim \frac{u^\top \Sigma_w u}{n} \chi_{K-1}^2 \left(\frac{nu^\top \tilde{\Sigma}_b^{(2)} u}{u^\top \Sigma_w u} \right),$$

where $\chi_r^2(\lambda)$ denotes a non-central chi-squared distribution with r degrees of freedom and non-centrality parameter λ . By Birgé (2001, Lemma 8.1), for every $\delta' \in (0, 1/2]$, we have with probability at least $1 - 2\delta'$ conditional on Y_1, \dots, Y_n , that

$$\begin{aligned} |u^\top (\hat{\Sigma}_b - \tilde{\Sigma}_b^{(2)}) u| &\leq \frac{u^\top \Sigma_w u}{n} \left\{ K + 2 \sqrt{\left(K + \frac{2nu^\top \tilde{\Sigma}_b^{(2)} u}{u^\top \Sigma_w u} \right) \log(1/\delta') + 2 \log(1/\delta')} \right\} \\ &\leq \frac{u^\top \Sigma_w u}{n} \{2K + 3 \log(1/\delta')\} + \sqrt{\frac{8u^\top \Sigma_w u u^\top \tilde{\Sigma}_b^{(2)} u \log(1/\delta')}{n}} \\ &\leq \frac{3(K + \log(1/\delta')) \|\Sigma_w\|_{\text{op}}}{n} + \sqrt{\frac{8 \|\Sigma_w\|_{\text{op}} \|\tilde{\Sigma}_b^{(2)}\|_{\text{op}} \log(1/\delta')}{n}}. \end{aligned}$$

Let \mathcal{N} be a $1/4$ -net of the sphere \mathbb{S}^{d-1} , which can be chosen to have cardinality at most 9^d (Vershynin, 2012, Lemma 5.2). Hence, through a union bound, and taking $\delta' := \delta/(8 \cdot 9^d)$, we have with probability at least $1 - \delta/4$ conditional on Y_1, \dots, Y_n that

$$\begin{aligned} \|\hat{\Sigma}_b - \tilde{\Sigma}_b^{(2)}\|_{\text{op}} &\leq 2 \max_{u \in \mathcal{N}} |u^\top (\hat{\Sigma}_b - \tilde{\Sigma}_b^{(2)}) u| \\ &\leq \frac{6R_2(K + \log(8 \cdot 9^d/\delta))}{n} + \sqrt{\frac{32R_1^2 R_2 \log(8 \cdot 9^d/\delta)}{n}}. \end{aligned} \quad (22)$$

Combining (19), (20), (21) and (22), we have that the desired bound on $\|\hat{\Sigma}_b - \Sigma_b\|_{\text{op}}$ occurs on an event with probability at least $1 - 3\delta/4$.

We now turn to control $\|\hat{\Sigma}_w - \Sigma_w\|_{\text{op}}$. Let $Q := \sum_{k \in [K]} (n_k^{-1} \mathbb{1}_{\{Y_i=k, Y_{i'}=k\}})_{i, i'=1}^n \in \mathbb{R}^{n \times n}$, so that

$$n\hat{\Sigma}_w = \sum_{i=1}^n Z_i Z_i^\top - \sum_{k=1}^K n_k \hat{\mu}_k \hat{\mu}_k^\top = Z^\top (I - Q) Z,$$

where $Z := (Z_1, \dots, Z_n)^\top$. It therefore follows again by Lemma 17 that $\hat{\Sigma}_w \mid Y_1, \dots, Y_n \sim n^{-1} \mathcal{W}_d(n - K, \Sigma_w)$. Another application of Muirhead (2009, Theorem 10.3.6) then yields for any $u \in \mathbb{S}^{d-1}$ that

$$u^\top \hat{\Sigma}_w u \mid Y_1, \dots, Y_n \sim \frac{u^\top \Sigma_w u}{n} \chi_{n-K}^2.$$

By Laurent and Massart (2000, Lemma 1), we have with probability at least $1 - 2\delta' = 1 - \delta/(4 \cdot 9^d)$ that

$$|u^\top (\hat{\Sigma}_w - \Sigma_w) u| \leq \frac{R_2}{n} \{K + 2\sqrt{n \log(1/\delta')} + 2 \log(1/\delta')\}.$$

Again, taking a union bound over the $1/4$ -net \mathcal{N} , we conclude that with probability at least $1 - \delta/4$, we have

$$\|\hat{\Sigma}_w - \Sigma_w\|_{\text{op}} \leq 2 \max_{u \in \mathcal{N}} |u^\top (\hat{\Sigma}_w - \Sigma_w) u| \leq \frac{4R_2}{n} \{K + \sqrt{n \log(8 \cdot 9^d/\delta)} + \log(8 \cdot 9^d/\delta)\},$$

as desired. \square

Proof of Theorem 3. Define $\delta := \binom{p}{d}^{-1} \varepsilon$. Since Algorithm 2 is permutation-equivariant, by a union bound, it suffices to show that for every $P \in \mathcal{P}_d$, with probability at least $1 - \delta$, the desired upper bound holds for $\|\psi((PX_i, Y_i)_{i \in [n]}) - (P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top\|_{\text{op}}$. Recall that $n_k := \sum_{i=1}^n \mathbb{1}_{\{Y_i=k\}}$. Write $\Sigma_{w,P} := P\Sigma_w P^\top$ and $\Sigma_{b,P} := P\Sigma_b P^\top$, $\hat{\mu}_{k,P} := n_k^{-1} \sum_{i: Y_i=k} PX_i$, $\hat{\mu}_P := n^{-1} \sum_{i=1}^n PX_i$ and

$$\hat{\Sigma}_{w,P} := \frac{1}{n} \sum_{i=1}^n (PX_i - \hat{\mu}_{Y_i,P})(PX_i - \hat{\mu}_{Y_i,P})^\top \quad \text{and} \quad \hat{\Sigma}_{b,P} := \sum_{k=1}^K \frac{n_k}{n} (\hat{\mu}_{k,P} - \hat{\mu}_P)(\hat{\mu}_{k,P} - \hat{\mu}_P)^\top.$$

Observe that since $n \geq Kd + 1$, we have $\max_{k \in [K]} n_k \geq d + 1$, so $\hat{\Sigma}_{w,P}$ is positive definite with probability 1. Thus, by the triangle inequality, for each $P \in \mathcal{P}_d$, we have

$$\begin{aligned} \|\psi((PX_i, Y_i)_{i \in [n]}) - (P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top\|_{\text{op}} &= \|\hat{\Sigma}_{w,P}^{-1} \hat{\Sigma}_{b,P} - \Sigma_{w,P}^{-1} \Sigma_{b,P}\|_{\text{op}} \\ &\leq \|\hat{\Sigma}_{w,P}^{-1} \hat{\Sigma}_{b,P} - \Sigma_{w,P}^{-1} \hat{\Sigma}_{b,P}\|_{\text{op}} + \|\Sigma_{w,P}^{-1} \hat{\Sigma}_{b,P} - \Sigma_{w,P}^{-1} \Sigma_{b,P}\|_{\text{op}}. \end{aligned} \quad (23)$$

By Proposition 8 and our hypothesis, there is an event Ω_P with probability at least $1 - \delta$, on which

$$\begin{aligned} \|\hat{\Sigma}_{w,P} - \Sigma_{w,P}\|_{\text{op}} &\leq \frac{4R_2 \{K + \log(8 \cdot 9^d/\delta)\}}{n} + 4R_2 \sqrt{\frac{\log(8 \cdot 9^d/\delta)}{n}} \leq \frac{1}{2R_2}. \\ \|\hat{\Sigma}_{b,P} - \Sigma_{b,P}\|_{\text{op}} &\lesssim_{R_1, R_2} \frac{K + d + \log(1/\delta)}{n} + \sqrt{\frac{d + \log(1/\delta)}{n}} \lesssim_{R_2} 1. \end{aligned} \quad (24)$$

Thus, for the first term in (23), by Weyl's inequality, on Ω_P , we have

$$\begin{aligned}
\|\hat{\Sigma}_{w,P}^{-1}\hat{\Sigma}_{b,P} - \Sigma_{w,P}^{-1}\hat{\Sigma}_{b,P}\|_{\text{op}} &\leq \|\Sigma_{w,P}^{-1} - \hat{\Sigma}_{w,P}^{-1}\|_{\text{op}}\|\hat{\Sigma}_{b,P}\|_{\text{op}} \\
&\leq \|\Sigma_{w,P}^{-1}\|_{\text{op}}\|\hat{\Sigma}_{w,P}^{-1}\|_{\text{op}}\|\hat{\Sigma}_{b,P}\|_{\text{op}}\|\hat{\Sigma}_{w,P} - \Sigma_{w,P}\|_{\text{op}} \\
&\leq \frac{(\|\Sigma_{b,P}\|_{\text{op}} + \|\hat{\Sigma}_{b,P} - \Sigma_{b,P}\|_{\text{op}})\|\hat{\Sigma}_{w,P} - \Sigma_{w,P}\|_{\text{op}}}{\lambda_{\min}(\Sigma_{w,P})(\lambda_{\min}(\Sigma_{w,P}) - \|\hat{\Sigma}_{w,P} - \Sigma_{w,P}\|_{\text{op}})} \\
&\leq \frac{R_2(R_1^2 + \|\hat{\Sigma}_{b,P} - \Sigma_{b,P}\|_{\text{op}})\|\hat{\Sigma}_{w,P} - \Sigma_{w,P}\|_{\text{op}}}{(1/R_2 - \|\hat{\Sigma}_{w,P} - \Sigma_{w,P}\|_{\text{op}})} \\
&\lesssim_{R_1, R_2} \|\hat{\Sigma}_{w,P} - \Sigma_{w,P}\|_{\text{op}} \lesssim_{R_2} \frac{K}{n} + \sqrt{\frac{d + \log(1/\delta)}{n}}, \quad (25)
\end{aligned}$$

where we used (24) in the penultimate inequality. For the second term in (23), we also have on Ω_P that

$$\|\Sigma_{w,P}^{-1}\hat{\Sigma}_{b,P} - \Sigma_{w,P}^{-1}\Sigma_{b,P}\|_{\text{op}} \leq \|\Sigma_{w,P}^{-1}\|_{\text{op}}\|\hat{\Sigma}_{b,P} - \Sigma_{b,P}\|_{\text{op}} \lesssim_{R_1, R_2} \frac{K}{n} + \sqrt{\frac{d + \log(1/\delta)}{n}}. \quad (26)$$

The desired result follows by combining (25) and (26), and using the fact that $\log(1/\delta) \leq d \log(ep/d) + \log(1/\varepsilon)$. \square

5.4 Proofs of Proposition 5 and Theorem 6

In the proof of Proposition 5, we show the convergence of the EM iterates $\hat{\mu}^{(t)}$ by analyzing their components parallel and orthogonal to μ^* separately. Writing $\eta := \mu^*/\|\mu^*\|$, let $\alpha_t \in \mathbb{R}$, $\beta_t \geq 0$ be defined by

$$\hat{\mu}^{(t)} = \alpha_t \eta + \beta_t \xi_t, \quad (27)$$

where $\xi_t \in \mathbb{S}^{d-1}$ is orthogonal to η . Our proof will combine several propositions that control α_t and β_t under different conditions. We begin by laying some groundwork and defining some quantities that will be used throughout this subsection.

First, it will be convenient to relabel the two classes as $\{-1, 1\}$ instead of $\{1, 2\}$. By the rotational symmetry of the problem, we may assume without loss of generality that $\mu^* = (s, 0, \dots, 0)^\top \in \mathbb{R}^d$ for some $s \geq 0$, and that the first n_L observations are labeled (i.e., $Y_i \neq 0$ for $i \in [n_L]$). We assume throughout this section that $s \leq r$ and $r \geq 1$. Let $\hat{\mu}_{n_L} := n_L^{-1} \sum_{i=1}^{n_L} Z_i Y_i$, with the convention that $\hat{\mu}_{n_L} := 0$ if $n_L = 0$, and define the function $f_{n_U} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ by

$$f_{n_U}(v) := \frac{1}{n_U} \sum_{i=n_L+1}^n Z_i \tanh\langle Z_i, v \rangle, \quad (28)$$

with $f_{n_U} := 0$ if $n_U = 0$. Throughout, and without further comment, we assume that $n = n_L + n_U \geq 2$. In this notation, the EM update (12) can be rewritten, defining the function $g_n : \mathbb{R}^d \rightarrow \mathbb{R}^d$, as

$$\hat{\mu}^{(t)} = g_n(\hat{\mu}^{(t-1)}) := \gamma \hat{\mu}_{n_L} + (1 - \gamma) f_{n_U}(\hat{\mu}^{(t-1)}).$$

The corresponding population quantities are

$$f(v) := \mathbb{E} Z_1 \tanh\langle v, Z_1 \rangle \quad \text{and} \quad g(v) := \gamma \mu^* + (1 - \gamma) f(v).$$

Writing $\Delta_{n_U} := f_{n_U} - f$, we have

$$g_n(v) = g(v) + (1 - \gamma)\Delta_{n_U}(v) + \gamma(\hat{\mu}_{n_L} - \mu^*). \quad (29)$$

For $\omega, \phi > 0$ and $r \geq 1$, we define the following two events that control the terms in the EM iteration involving the unlabeled and labeled data respectively:

$$\Omega_1(\omega) := \left\{ \sup_{v \in \mathbb{R}^d} \|g_n(v)\| \leq 2(r + \sqrt{d}) \right\} \cap \left\{ \sup_{\substack{\|v\| \leq 2(r + \sqrt{d}) \\ v \neq 0}} \frac{\|\Delta_{n_U}(v)\|}{\|v\|} \leq \omega \right\} \quad (30)$$

$$\Omega_2(\phi) := \{ \|\hat{\mu}_{n_L} - \mu^*\| \leq \phi \}.$$

Proposition 9. *There exists $C_r > 0$, depending only on r , such that for any $\delta \in (2e^{-n}, 1]$ and $\omega = C_r \sqrt{\frac{d \log n + \log(1/\delta)}{n_U}}$, we have $\mathbb{P}(\Omega_1(\omega)^c) \leq \delta$. Moreover, for any $\delta \in (0, 1]$ and for $\phi = \sqrt{\frac{2d + 3 \log(1/\delta)}{n_L}}$, we have $\mathbb{P}(\Omega_2(\phi)^c) \leq \delta$.*

Proof. For any $v \in \mathbb{R}^d$,

$$\begin{aligned} \|g_n(v)\| &= \|(1 - \gamma)f_{n_U}(v) + \gamma\hat{\mu}_{n_L}\| \leq (1 - \gamma) \cdot \frac{1}{n_U} \sum_{i=n_L+1}^n \|Z_i\| + \frac{\gamma}{n_L} \sum_{i=1}^{n_L} \|Z_i\| \\ &= \frac{1}{n} \sum_{i=1}^n \|Z_i\| \leq \left(\frac{1}{n} \sum_{i=1}^n \|Z_i\|^2 \right)^{1/2}. \end{aligned}$$

Since $\sum_{i=1}^n \|Z_i\|^2 \sim \chi_{nd}^2(ns^2)$, by [Birgé \(2001, Lemma 8.1\)](#), we have with probability at least $1 - \delta/2$ that

$$\begin{aligned} \sup_{v \in \mathbb{R}^d} \|g_n(v)\|^2 &\leq d + s^2 + 2\sqrt{\frac{(d + 2s^2) \log(2/\delta)}{n}} + \frac{2 \log(2/\delta)}{n} \\ &\leq 2d + 3s^2 + \frac{3 \log(2/\delta)}{n} \leq 4(r + \sqrt{d})^2. \end{aligned} \quad (31)$$

Also, by a very similar argument as in the proof of [Wu and Zhou \(2022, Theorem 4\)](#), we have with probability at least $1 - \delta/2$ that

$$\sup_{\substack{\|v\| \leq 2(r + \sqrt{d}) \\ v \neq 0}} \frac{\|\Delta_{n_U}(v)\|}{\|v\|} \leq C_r \sqrt{\frac{d \log n + \log(1/\delta)}{n_U}}, \quad (32)$$

for some $C_r > 0$ depending only on r . The first claim follows by combining (31) and (32).

For the second claim, we have $\hat{\mu}_{n_L} \sim N_d(\mu^*, n_L^{-1}I_d)$. Hence, by [Laurent and Massart \(2000, Lemma 1\)](#), we have

$$\begin{aligned} \mathbb{P}(\Omega_2(\phi)^c) &= \mathbb{P}(n_L \|\hat{\mu}_{n_L} - \mu^*\|^2 > n_L \phi^2) \\ &\leq \mathbb{P}(n_L \|\hat{\mu}_{n_L} - \mu^*\|^2 > d + 2\sqrt{d \log(1/\delta)} + 2 \log(1/\delta)) \leq \delta, \end{aligned}$$

as required. \square

For any $a \in \mathbb{R}$, $b \in [0, \infty)$ and $\xi \in \mathbb{S}^{d-1}$ that is orthogonal to η , we define $F(a, b) := \eta^\top f(a\eta + b\xi)$ and $G(a, b) := \|(I_d - \eta\eta^\top)f(a\eta + b\xi)\|$. Note that the distribution of Z_1 is orthogonally invariant along the axis μ^* ; in other words, if $P \in \mathbb{R}^{d \times d}$ is orthogonal and has μ^* as an eigenvector with eigenvalue 1, then $PZ_1 \stackrel{d}{=} Z_1$. It follows that $f(a\eta + b\xi)$, and hence $F(a, b)$ and $G(a, b)$, do not depend on ξ . We remark that

$$f(\alpha_t\eta + \beta_t\xi_t) = F(\alpha_t, \beta_t)\eta + G(\alpha_t, \beta_t)\xi'_{t+1}$$

for some $\xi'_{t+1} \in \mathbb{S}^{d-1}$ that is orthogonal to η .

Proposition 10 controls the magnitude of the component β_t of the EM algorithm iterates that is orthogonal to the signal direction η . We define $\zeta := \omega\gamma^{-1/2} \wedge \omega^{1/2}$.

Proposition 10. *Assume that $\phi\gamma^{1/2} \leq \omega \leq \min\{1/12, 1/(r+3)\}$ and that $\|\hat{\mu}^{(0)}\| \leq r+3$. On the event $\Omega_1(\omega) \cap \Omega_2(\phi)$, we have*

$$\limsup_{t \rightarrow \infty} \beta_t \leq 60(\zeta \vee r\omega).$$

Moreover, on the same event, if $\beta_{t_0} \leq 60(\zeta \vee r\omega)$ for some $t_0 \in \mathbb{N}_0$, then $\beta_t \leq 60(\zeta \vee r\omega)$ for all $t \geq t_0$.

Proof. We first claim that on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$, we have $\|\hat{\mu}^{(t)}\| \leq r+3$ for all $t \in \mathbb{N}_0$. The case $t=0$ is true by the assumption on the initializer $\hat{\mu}^{(0)}$, and if the claim holds for $t \in \mathbb{N}_0$, then since $2(r + \sqrt{d}) \geq 2(r+1) \geq r+3$, we have on $\Omega_1(\omega) \cap \Omega_2(\phi)$ that

$$\begin{aligned} \|\hat{\mu}^{(t+1)}\| &\leq (1-\gamma)\{|F(\alpha_t, \beta_t)| + G(\alpha_t, \beta_t) + \|\Delta_{n_U}(\hat{\mu}^{(t)})\|\} + \gamma(s + \phi) \\ &\leq s + 2\sqrt{2/\pi} + \omega\|\hat{\mu}^{(t)}\| + \gamma\phi \leq r + 2 + \frac{\|\hat{\mu}^{(t)}\|}{r+3} \leq r + 3, \end{aligned}$$

where the second inequality uses Wu and Zhou (2022, Lemma 5(5)). Moreover, from (29), we have on $\Omega_1(\omega) \cap \Omega_2(\phi)$ that for $t \in \mathbb{N}$,

$$\begin{aligned} \beta_{t+1} &= \|(I_d - \eta\eta^\top)\{(1-\gamma)(f(\hat{\mu}^{(t)}) + \Delta_{n_U}(\hat{\mu}^{(t)})) + \gamma\hat{\mu}_{n_L}\}\| \\ &\leq (1-\gamma)\{G(\alpha_t, \beta_t) + \omega(|\alpha_t| + \beta_t)\} + \gamma\phi \\ &\leq \beta_t(1-\gamma)\left\{1 + \omega - \frac{(\alpha_t^2 + \beta_t^2) \wedge 1}{6}\right\} + \gamma\phi + \omega|\alpha_t|, \end{aligned} \tag{33}$$

where the final bound uses Wu and Zhou (2022, Lemma 5(8)). If $\alpha_t^2 + \beta_t^2 > 1$ or $\gamma > 1/2$, then using the fact that $\omega \leq 1/12$, we have from (33) that

$$\beta_{t+1} \leq \frac{11}{12}\beta_t + \gamma\phi + (r+3)\omega \leq \frac{11}{12}\beta_t + (r+4)\omega. \tag{34}$$

On the other hand, if $\alpha_t^2 + \beta_t^2 \leq 1$ and $\gamma \leq 1/2$, then

$$\beta_{t+1} \leq \beta_t \left(1 + \omega - \gamma - \frac{\alpha_t^2 + \beta_t^2}{12}\right) + \gamma\phi + \omega|\alpha_t|. \tag{35}$$

Note that the right-hand side of (34) is increasing in β_t and the right-hand side of (35) is increasing in β_t for $\alpha_t^2 + \beta_t^2 \leq 1$ and $\gamma \leq 1/2$. Combining (34) and (35), denoting $\beta_\infty :=$

$\limsup_{t \rightarrow \infty} \beta_t$ and using the fact that $0 \leq \frac{3}{\beta_\infty}(\omega - |\alpha_t| \beta_\infty / 6)^2 = 3\omega^2 / \beta_\infty - \omega |\alpha_t| + \alpha_t^2 \beta_\infty / 12$, we have

$$\beta_\infty \leq \max \left\{ \frac{11}{12} \beta_\infty + (r+4)\omega, \beta_\infty \left(1 + \omega - \gamma - \frac{\beta_\infty^2}{12} \right) + \gamma\phi + \frac{3\omega^2}{\beta_\infty} \right\}. \quad (36)$$

From the first term in the maximum in (36), we obtain

$$\beta_\infty \leq (1r + 38)\omega \leq 60r\omega. \quad (37)$$

From the second term in the maximum in (36), we obtain

$$\beta_\infty \left(\gamma - \omega + \frac{\beta_\infty^2}{12} \right) \leq \gamma\phi + \frac{3\omega^2}{\beta_\infty}. \quad (38)$$

If $\gamma < 2\omega$, then from (38),

$$\beta_\infty \leq 5\omega^{1/2} \leq 5\sqrt{2}\zeta, \quad (39)$$

since otherwise we would have that the left-hand side would be at least $(-5 + 125/12)\omega^{3/2}$ and the right-hand side would at most $(\sqrt{2} + 3/5)\omega^{3/2}$, contradicting the inequality. On the other hand, if $\gamma \geq 2\omega$, then we derive from (38) that

$$\beta_\infty^2 - 2\phi\beta_\infty - \frac{6\omega^2}{\gamma} \leq 0.$$

Solving this inequality, we find that

$$\beta_\infty \leq \phi + \sqrt{\phi^2 + 6\omega^2/\gamma} \leq 4\omega\gamma^{-1/2} = 4\zeta. \quad (40)$$

The first claim of the proposition follows by combining (36), (37), (39) and (40). We now prove the second claim by induction on t . The base case $t = t_0$ is true by assumption, so we assume that $\beta_t \leq 60(\zeta \vee r\omega)$ for some $t \geq t_0$. Again we consider two cases. If $\beta_t \leq 36(\zeta \vee r\omega)$, then from (33) and using that $|\alpha_t| \leq \|\hat{\mu}^{(t)}\| \leq r + 3$ for $t \geq 2$,

$$\beta_{t+1} \leq \beta_t(1 + \omega) + \gamma\phi + \omega|\alpha_t| \leq 44(\zeta \vee r\omega),$$

as desired. On the other hand, if $\beta_t > 36(\zeta \vee r\omega)$, then combining (34) and (35), we obtain that

$$\beta_{t+1} \leq \max \left\{ 60(\zeta \vee r\omega), \beta_t \left(1 + \omega - \gamma - \frac{\beta_t^2}{12} \right) + \gamma\phi + \frac{3\omega^2}{\beta_t} \right\}.$$

It suffices to show the second term in the maximum is no larger than β_t . To this end, if $\gamma \leq 2\omega$, then $\zeta^2 \leq \omega \leq 2\zeta^2$, and so

$$\begin{aligned} \beta_t \left(1 + \omega - \gamma - \frac{\beta_t^2}{12} \right) + \gamma\phi + \frac{3\omega^2}{\beta_t} &\leq \beta_t + 2\zeta^2\beta_t - \frac{\beta_t^3}{12} + \gamma^{1/2}\omega + \frac{3\omega^2}{\beta_t} \\ &\leq \beta_t + 120(\zeta \vee r\omega)\zeta^2 - \frac{36^3(\zeta \vee r\omega)^3}{12} + 4\zeta^3 + \frac{\zeta^3}{3} \leq \beta_t. \end{aligned}$$

On the other hand, if $\gamma > 2\omega$, then $\zeta = \omega\gamma^{-1/2} \geq \phi$, and so

$$\begin{aligned} \beta_t \left(1 + \omega - \gamma - \frac{\beta_t^2}{12} \right) + \gamma\phi + \frac{3\omega^2}{\beta_t} &\leq \beta_t - \frac{\gamma\beta_t}{2} + \gamma\phi + \frac{3\omega^2}{\beta_t} \\ &\leq \beta_t - 18\gamma\zeta + \gamma\zeta + \frac{\omega^2}{12\zeta} \leq \beta_t, \end{aligned}$$

as desired, which completes the induction. \square

The following result bounds the magnitude of the signal component, α_t , of the EM iterates.

Proposition 11. *Assume that $\phi\gamma^{1/2} \leq \omega \leq \min\{1/12, 1/(r+3)\}$ and that $\|\hat{\mu}^{(0)}\| \leq r+3$. Then there exists $C_r > 0$, depending only on r , such that on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$, we have*

$$\limsup_{t \rightarrow \infty} |\alpha_t| \leq C_r(\zeta \vee s).$$

Proof. By definition of α_{t+1} and (29), we have for every $t \in \mathbb{N}_0$ that

$$\alpha_{t+1} = \eta^\top \left\{ (1 - \gamma)(f(\hat{\mu}^{(t)}) + \Delta_{n_U}(\hat{\mu}^{(t)})) + \gamma \hat{\mu}_{n_L} \right\}.$$

Thus, by the first claim in the proof of Proposition 10, we have on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$ that

$$|\alpha_{t+1} - (1 - \gamma)F(\alpha_t, \beta_t) - \gamma s| \leq (1 - \gamma)\omega(|\alpha_t| + \beta_t) + \gamma\phi \quad (41)$$

for every $t \in \mathbb{N}_0$. From Wu and Zhou (2022, Lemma 5(1) and Lemma 5(7)), $\alpha \mapsto F(\alpha, \beta)$ is an increasing and odd function satisfying $|F(\alpha, \beta) - F(\alpha, 0)| \leq (1 + s^2)|\alpha|\beta^2$ for every $\alpha, \beta \in \mathbb{R}$. Hence, by (41), we have on $\Omega_1(\omega) \cap \Omega_2(\phi)$ that

$$|\alpha_{t+1}| \leq (1 - \gamma) \left\{ F(|\alpha_t|, 0) + (1 + s^2)|\alpha_t|\beta_t^2 + \omega(|\alpha_t| + \beta_t) \right\} + \gamma(s + \phi). \quad (42)$$

Note the right-hand side of (42) is increasing in $|\alpha_t|$. Define $\alpha_\infty := \limsup_{t \rightarrow \infty} |\alpha_t|$, so that $\alpha_\infty \leq r + 3$ on $\Omega_1(\omega) \cap \Omega_2(\phi)$, again by the first claim in the proof of Proposition 10. We may also assume that $\alpha_\infty > s$, because otherwise the result is clear. Since $\alpha \mapsto F(\alpha, 0)$ is continuous, we have from (42) that on $\Omega_1(\omega) \cap \Omega_2(\phi)$,

$$\alpha_\infty \leq (1 - \gamma) \left\{ F(\alpha_\infty, 0) + (1 + r^2)\alpha_\infty\beta_\infty^2 + \omega(\alpha_\infty + \beta_\infty) \right\} + \gamma(s + \phi), \quad (43)$$

where we recall that $\beta_\infty := \limsup_{t \rightarrow \infty} \beta_t$. Define $q : [0, \infty) \rightarrow \mathbb{R}$ by

$$q(\alpha) := \begin{cases} F(\alpha, 0)/\alpha & \text{if } \alpha \neq 0 \\ 1 + s^2 & \text{if } \alpha = 0. \end{cases} \quad (44)$$

By Lemma 20, we have $q(s) = 1$ (which confirms that μ_* is a fixed point of the population EM iteration), and that $q'(\alpha) \leq -c_r\alpha$ for all $\alpha \in (0, r]$, where $c_r \in (0, 1]$ depends only on r . Thus, dividing both sides of (43) by α_∞ , we have

$$\begin{aligned} 1 &\leq (1 - \gamma) \left\{ q(s) + \int_s^{\alpha_\infty} q'(\alpha) d\alpha + (1 + r^2)\beta_\infty^2 + \omega \left(1 + \frac{\beta_\infty}{\alpha_\infty} \right) \right\} + \frac{\gamma(s + \phi)}{\alpha_\infty} \\ &\leq (1 - \gamma) \left\{ 1 - \frac{c_r}{2}(\alpha_\infty^2 - s^2) + (1 + r^2)\beta_\infty^2 + \omega \left(1 + \frac{\beta_\infty}{\alpha_\infty} \right) \right\} + \frac{\gamma(s + \phi)}{\alpha_\infty}. \end{aligned} \quad (45)$$

Now $\beta_\infty \leq 60(1+r)\zeta$ by Proposition 10. We now claim that $\alpha_\infty \leq 4s + 120c_r^{-1/2}(1+r)^2\zeta$. Indeed, assuming the contrary, we would have $c_r\alpha_\infty^2/4 > c_r s^2/2 + (1+r^2)\beta_\infty^2$ and $\beta_\infty/\alpha_\infty < 1$. Hence from (45), we have

$$1 \leq (1 - \gamma) \left(1 - \frac{c_r\alpha_\infty^2}{4} \right) + \frac{\gamma s}{\alpha_\infty} + \left(2 + \frac{\gamma^{1/2}}{\alpha_\infty} \right) \omega.$$

We consider two cases. First, if $\gamma \leq 4\omega$, then $2\zeta \geq \omega^{1/2} \geq \gamma^{1/2}/2$ and hence

$$\begin{aligned} & (1 - \gamma) \left(1 - \frac{c_r \alpha_\infty^2}{4} \right) + \frac{\gamma s}{\alpha_\infty} + \left(2 + \frac{\gamma^{1/2}}{\alpha_\infty} \right) \omega \\ & \leq \max \left\{ 1 - \frac{c_r \alpha_\infty^2}{4}, 0 \right\} + \frac{\gamma}{4} + 3\omega \leq \max \left\{ 1 - \zeta^2, \frac{\gamma}{4} + 3\omega \right\} < 1, \end{aligned}$$

a contradiction. Second, if $\gamma > 4\omega$, then $\zeta = \omega\gamma^{-1/2}$ and

$$(1 - \gamma) \left(1 - \frac{c_r \alpha_\infty^2}{4} \right) + \frac{\gamma s}{\alpha_\infty} + \left(2 + \frac{\gamma^{1/2}}{\alpha_\infty} \right) \omega \leq 1 - \gamma + \frac{\gamma}{4} + 2\omega + \frac{\gamma\zeta}{\alpha_\infty} < 1,$$

again a contradiction. This establishes the claimed upper bound on α_∞ . \square

Recall the definitions $L(\mu, \mu^*) = \|\mu - \mu^*\| \wedge \|\mu + \mu^*\|$. Our next result shows that if α_t ever becomes sufficiently large, then improved bounds can be derived on the limiting behaviour of α_t , β_t and $L(\hat{\mu}^{(t)}, \mu^*)$.

Proposition 12. *Assume that $\gamma \in [0, 1/2)$. Given any $c > 0$, there exists $C, c_1 > 0$, depending only on r and c , such that if $|\alpha_{t_0}| \geq cs$, $\beta_{t_0} \leq 60(\zeta \vee r\omega)$ for some iteration t_0 , and $\phi\gamma^{1/2} \leq \omega \leq c_1$ and $s \geq C\zeta$, then on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$, we have*

$$\limsup_{t \rightarrow \infty} |\alpha_t - s| \lesssim_{r,c} \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}}, \quad (46)$$

$$\limsup_{t \rightarrow \infty} \beta_t \lesssim_{r,c} \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}}, \quad (47)$$

$$\limsup_{t \rightarrow \infty} L(\hat{\mu}^{(t)}, \mu^*) \lesssim_{r,c} \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}}. \quad (48)$$

Proof. By flipping the sign of μ^* if necessary, we may assume without loss of generality that $\alpha_{t_0} \geq 0$ and that $c_1 \leq \min\{1/12, 1/(r+3)\}$. From (41) and the argument immediately below it, we have

$$\begin{aligned} & (1 - \gamma) \left\{ F(\alpha_t, 0) - (1 + s^2)\alpha_t\beta_t^2 - \omega(\alpha_t + \beta_t) \right\} + \gamma(s - \phi) \leq \alpha_{t+1} \\ & \leq (1 - \gamma) \left\{ F(\alpha_t, 0) + (1 + s^2)\alpha_t\beta_t^2 + \omega(\alpha_t + \beta_t) \right\} + \gamma(s + \phi). \end{aligned} \quad (49)$$

For any t such that $\alpha_t \geq cs$, since $\beta_t \leq 60r\omega^{1/2}$ by Proposition 10, we have that

$$(1 + s^2)\alpha_t\beta_t^2 + \omega(\alpha_t + \beta_t) \leq \left\{ (1 + r^2)60^2r^2\omega + \omega \left(1 + \frac{60s\zeta}{cs} \right) \right\} \alpha_t \leq c''\omega\alpha \quad (50)$$

where $c'' := 60^2(1 + r^2)r^2 + (1 + 60rc^{-1}C^{-1})$. Moreover, if $\alpha_t \geq cs$, then

$$\gamma\phi \leq \omega\gamma^{1/2} \leq \begin{cases} \gamma\zeta \leq C^{-1}\gamma s & \text{if } \gamma \geq \omega \\ \omega^{3/2} = \omega\zeta \leq c^{-1}C^{-1}\omega\alpha_t & \text{otherwise.} \end{cases} \quad (51)$$

Let $c' := c'' + 2c^{-1}C^{-1}$ and define functions $H, L : [0, \infty) \rightarrow \mathbb{R}$ by

$$\begin{aligned} H(\alpha) &:= (1 - \gamma) \left\{ F(\alpha, 0) + c'\omega\alpha \right\} + \gamma(s + \phi), \\ L(\alpha) &:= (1 - \gamma) \left\{ F(\alpha, 0) - c'\omega\alpha \right\} + \gamma \max(s - \phi, s/2), \end{aligned}$$

From (49), (50) and (51), we obtain that for $\alpha_t \geq cs$ and $C \geq 2$,

$$L(\alpha_t) \leq \alpha_{t+1} \leq H(\alpha_t). \quad (52)$$

Define auxiliary sequences $(\alpha_t^+)_{t \geq t_0}$ and $(\alpha_t^-)_{t \geq 0}$ by $\alpha_{t_0}^+ := \alpha_{t_0} =: \alpha_{t_0}^-$ and for $t \geq t_0$,

$$\alpha_{t+1}^+ := H(\alpha_t^+) \quad \text{and} \quad \alpha_{t+1}^- := L(\alpha_t^-).$$

We first derive some properties of the two recursion maps H and L . For the former, we have by Wu and Zhou (2022, Lemma 3) that F , and hence H , is increasing and concave on $[0, \infty)$ with $H(0) > 0$ when $\gamma > 0$ and $H'(0) > \partial_1 F(0, 0) > 1$ when $\gamma = 0$. Moreover, since F is bounded, we can choose $c_1 > 0$, depending only on r and c , such that $\lim_{\alpha \rightarrow \infty} H'(\alpha) = (1 - \gamma)c'\omega \leq (1 - \gamma)c'c_1 < 1/2$. On the other hand, we have $L(0) > 0$ when $\gamma > 0$. When $\gamma = 0$, we have $\omega^{1/2} = \zeta \leq s/C$, which means that after increasing $C \equiv C(r, c) > 0$ if necessary, $L'(0) = \partial_1 F(0, 0) - c'\omega \geq 1 + s^2 - c's^2/C^2 > 1$. By Wu and Zhou (2022, Lemma 3), $\alpha \mapsto F(\alpha, 0)$ is differentiable, increasing and concave for $\alpha \in [0, \infty)$. Reducing $c_1 \equiv c_1(r, c) > 0$ if necessary to ensure that $c_1 \leq \partial_1 F(r + 3, 0)/c'$, we have for $\alpha \in [0, r + 3]$ that

$$L'(\alpha) = \partial_1 F(\alpha, 0) - c'\omega \geq \partial_1 F(r + 3, 0) - c'c_1 \geq 0.$$

In other words, L is increasing on $[0, r + 3]$, and moreover, similarly to H , it is also concave on this interval. Finally, we claim that for $\tilde{c} := \min\{c, 32(3 + r^4)/3\}^{-1/2}$, and $\tilde{\alpha} := \tilde{c}s$, we have $L(\tilde{\alpha}) \geq \tilde{\alpha}$. To verify this, we note by Wu and Zhou (2022, Lemma 3(3) and Equation (98)), we have

$$L(\tilde{\alpha}) \geq (1 - \gamma)\tilde{\alpha} \left\{ 1 + s^2 - \frac{8}{3}(3 + r^4)\tilde{\alpha}^2 - c'\omega \right\} + 2\gamma\tilde{\alpha} \geq (1 - \gamma)\tilde{\alpha} \left(1 + \frac{3s^2}{4} - c'\omega \right) + 2\gamma\tilde{\alpha}. \quad (53)$$

To control the right-hand side of (53), if $\gamma \leq c'\omega$, we have $\zeta = \omega^{1/2} \wedge \omega\gamma^{-1/2} \geq (\omega/c')^{1/2}$. Hence, from the condition $s \geq C\zeta$, if we choose $C > 2c'$ (which is possible because c' is a decreasing function of C), we have $c'\omega \leq (c's/C)^2 \leq s^2/4$ and consequently the right-hand side of (53) is at least $\tilde{\alpha}$. If $\gamma > c'\omega$, then we have $L(\tilde{\alpha}) \geq (1 - \gamma)^2\tilde{\alpha} + 2\gamma\tilde{\alpha} \geq \tilde{\alpha}$ as desired. This establishes the claim.

We now show by induction that for all $t \geq t_0$,

$$\tilde{\alpha} \leq \alpha_t^- \leq \alpha_t \leq \alpha_t^+. \quad (54)$$

The base case is clear by the definition of $\alpha_{t_0}^-$ and $\alpha_{t_0}^+$ above. Now suppose that (54) holds for some iteration $t \geq 0$, so in particular, (52) applies.

Using the monotonicity of H and (52), we have $\alpha_{t+1} \leq H(\alpha_t) \leq H(\alpha_t^+) = \alpha_{t+1}^+$. Observe that $\alpha_t \leq \|\hat{\mu}^{(t)}\| \leq r + 3$ by the proof of Proposition 10. Using the monotonicity of L on $[0, r + 3]$, we find that $\alpha_{t+1} \geq L(\alpha_t) \geq L(\alpha_t^-) = \alpha_{t+1}^-$. Moreover,

$$\alpha_{t+1}^- = L(\alpha_t^-) \geq L(\tilde{\alpha}) \geq \tilde{\alpha},$$

which completes the induction.

To prove (46), we will analyze the sequences $(\alpha_t^+)_{t \geq t_0}$ and $(\alpha_t^-)_{t \geq t_0}$, which sandwich $(\alpha_t)_{t \geq t_0}$. We start by considering the behaviour of $(\alpha_t^+)_{t \geq t_0}$. The properties of H derived above mean that we can apply Lemma 19 to obtain that α_t^+ converges to a limit, denoted α^+ , satisfying $\alpha^+ = H(\alpha^+)$. By Lemma 20, we have $F(s, 0) = s$ and so $H(s) = (1 -$

$\gamma)(s + c'\omega s) + \gamma(s + \phi) > s$. Hence from Lemma 19 again, we have $\alpha^+ > s$. On the other hand, since $F(\alpha, 0) \leq \mathbb{E}|Z_{1,1}| \leq (\mathbb{E}Z_{1,1}^2)^{1/2} \leq (1 + s^2)^{1/2} \leq 1 + r$, we have $\alpha^+ = H(\alpha^+) \leq (1 - \gamma)(1 + r) + \alpha^+/2 + \gamma r + 1 \leq r + 2 + \alpha^+/2$, so $\alpha^+ \leq 2r + 4$. Recalling the definition of q from (44), by Lemma 20 again, we have

$$q(\alpha^+) = q(s) + \int_s^{\alpha^+} q'(\alpha) d\alpha \leq 1 - c_2((\alpha^+)^2 - s^2)$$

for some $c_2 > 0$ depending only on r . Consequently,

$$\begin{aligned} \alpha^+ &= H(\alpha^+) = (1 - \gamma)\alpha^+ \{q(\alpha^+) + c'\omega\} + \gamma(s + \phi) \\ &\leq (1 - \gamma)\alpha^+ \{1 - c_2((\alpha^+)^2 - s^2) + c'\omega\} + \gamma(s + \omega\gamma^{-1/2}), \end{aligned}$$

so

$$(\alpha^+)^2 - s^2 \leq \frac{c'\omega}{c_2} - \frac{\gamma}{(1 - \gamma)c_2} \frac{\alpha^+ - s - \omega\gamma^{-1/2}}{\alpha^+}. \quad (55)$$

We now prove that

$$\alpha^+ - s \lesssim_{r,c} \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}} \quad (56)$$

by considering two cases. If $\alpha^+ \leq 2s$, then from (55), we have

$$\alpha^+ - s \leq \frac{(1 - \gamma)c'\omega + \gamma^{1/2}\omega/\alpha^+}{(1 - \gamma)c_2(\alpha^+ + s) + \gamma/\alpha^+} \lesssim_{r,c} \frac{\omega}{s} \cdot \frac{1 + \gamma^{1/2}/s}{1 + \gamma/s^2} \lesssim \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}}.$$

On the other hand, if $\alpha^+ > 2s$, then we have from (55) again that

$$\frac{3(\alpha^+)^2}{4} + \frac{\gamma}{2c_2} \leq (\alpha^+)^2 - s^2 + \frac{\gamma(\alpha^+ - s)}{(1 - \gamma)c_2\alpha^+} \leq \frac{(c' + 2\gamma^{1/2}/\alpha^+)\omega}{c_2}. \quad (57)$$

In particular, $\gamma/2 \leq (c' + \gamma^{1/2}/s)\omega \leq c'\omega + C^{-1}(\gamma^{1/2}\omega^{1/2} \vee \gamma)$, so $\gamma \lesssim_{r,c} \omega$ and $\alpha^+ \gtrsim_{r,c} \zeta \gtrsim_{r,c} \gamma^{1/2}$. Consequently from (57),

$$\alpha^+ - s \leq \alpha^+ \lesssim_{r,c} \omega^{1/2} \lesssim_{r,c} \frac{\omega}{\alpha^+} \wedge \frac{\omega}{\gamma^{1/2}} \lesssim \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}},$$

which establishes (56). We now consider $(\alpha_t^-)_{t \geq t_0}$. Define $\tilde{L} : [0, \infty) \rightarrow [0, \infty)$ by $\tilde{L}(\alpha) := L(\alpha \wedge (r + 3))$. Since $\alpha_t^- \leq \alpha_t \leq r + 3$ for all $t \geq 0$, we have $\alpha_{t+1}^- = \tilde{L}(\alpha_t^-)$ for all $t \geq t_0$. From the properties of L derived above, we see that \tilde{L} satisfies the conditions of Lemma 19, and hence α_t^- converges to a limit, denoted α^- , satisfying $\alpha^- = \tilde{L}(\alpha^-) = L(\alpha^-)$. By Lemma 20, $F(s, 0) = s$, so we have $\tilde{L}(s) = L(s) \leq (1 - \gamma)(s - c'\omega s) + \gamma s < s$, so by Lemma 19, we must have $\alpha^- < s$. By Lemma 20 again, we have

$$q(\alpha^-) = q(s) - \int_{\alpha^-}^s q'(\alpha) d\alpha \geq 1 + c'_2(s^2 - (\alpha^-)^2),$$

where $c'_2 > 0$ depends only on r . Consequently, we have

$$\alpha^- = L(\alpha^-) \geq (1 - \gamma)\alpha^- \{1 + c'_2(s^2 - (\alpha^-)^2) - c'\omega\} + \gamma(s - \omega\gamma^{-1/2}),$$

which after rearranging and using the fact that $\alpha^- \geq \tilde{\alpha} \gtrsim_r s$ leads to

$$s - \alpha^- \leq \frac{(1 - \gamma)c'\omega + \gamma^{1/2}\omega/\alpha^-}{(1 - \gamma)c'_2(s + \alpha^-) + \gamma/\alpha^-} \lesssim_{r,c} \frac{\omega}{s} \cdot \frac{1 + \gamma^{1/2}/s}{1 + \gamma/s^2} \lesssim \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}}. \quad (58)$$

Combining (54), (56) and (58), we have established (46).

We now turn to prove (47). By increasing C if necessary, we have for all sufficiently large t that $s/2 \leq \alpha_t \leq 2s$. Consequently, we have by (33) that for all large t ,

$$\beta_{t+1} \leq \beta_t(1-\gamma) \left(1 + \omega - \frac{s^2/4 \wedge 1}{6} \right) + \gamma\phi + 2s\omega \leq \beta_t(1-\gamma)(1+\omega - c_3s^2) + (\gamma^{1/2} + 2s)\omega, \quad (59)$$

with $c_3 := 1/(6r^2 + 24)$. Denote $\beta_\infty := \limsup_{t \rightarrow \infty} \beta_t$. If $\gamma \leq 2\omega$, then $\zeta \in [(\omega/2)^{1/2}, \omega^{1/2}]$ so $\omega \leq 2s^2/C^2 \leq c_3s^2/2$ for C sufficiently large. Hence, from (59), and the fact that $\gamma^{1/2} \leq 2\zeta \leq s$,

$$\beta_\infty \leq \frac{(\gamma^{1/2} + 2s)\omega}{c_3s^2/2} \lesssim_r \frac{\omega}{s} = \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}}.$$

On the other hand, if $\gamma > 2\omega$, then $(1-\gamma)(1+\omega - c_3s^2) \leq 1 - \gamma/2 - c_3s^2/2$ and from (59), we obtain

$$\beta_\infty \leq \frac{(\gamma^{1/2} + 2s)\omega}{(\gamma + c_3s^2)/2} \lesssim_r \frac{\omega}{s} \cdot \frac{1 + \gamma^{1/2}/s}{1 + \gamma/s^2} \lesssim \frac{\omega}{s} \wedge \frac{\omega}{\gamma^{1/2}}.$$

Combining the above two cases establishes (47).

Finally, recalling the decomposition of $\hat{\mu}^{(t)}$ in (27), we see that (48) follows immediately from (46) and (47). \square

Next, we show that provided the initialization is not too uncorrelated with the true parameter, $|\alpha_t|$ reaches a level that makes Proposition 12 applicable after a sufficient number of iterations.

Proposition 13. *Assume that $n \geq 3$, that $\phi\gamma^{1/2} \leq \omega \leq \min\{1/12, 1/(r+3)\}$ and that $\gamma \in [0, 1/2)$. Suppose that $\hat{\mu}^{(0)}$ is chosen such that $c'(\zeta \vee r\omega) \leq \|\hat{\mu}^{(0)}\| \leq 60(\zeta \vee r\omega)$ for some $c' \in (0, 1)$ and that $|\langle \hat{\mu}^{(0)} / \|\hat{\mu}^{(0)}\|, \eta \rangle| \geq \sqrt{1/(d \log n_U)}$. Then there exist $c_4, c_5 > 0$, depending only on r and c' , such that if $s \geq c_4\zeta\sqrt{d \log n_U}$, then on $\Omega_1(\omega) \cap \Omega_2(\phi)$, we have*

$$|\alpha_t| \geq c_5s$$

for some $t > 0$.

Proof. By flipping the sign of μ^* if necessary, we may assume without loss of generality that $\alpha_0 \geq 0$. Assuming that the desired result is not true, we will prove by induction that on $\Omega_1(\omega) \cap \Omega_2(\phi)$, (a) $\alpha_t/\beta_t \geq c'/(60\sqrt{d \log n_U})$ and (b) $\alpha_{t+1} \geq (1 + \omega\sqrt{d \log n_U})\alpha_t$ for all $t \geq 0$. We show this by first verifying the base case of (a), then proving that (a) implies (b) for each t , and finally proving that $\alpha_{t+1}/\beta_{t+1} \geq 1/(60\sqrt{d \log n_U})$ once (b) holds for a given t .

For the base case, from the assumption on $\hat{\mu}^{(0)}$, we have

$$\frac{\alpha_0}{\beta_0} \geq \frac{\alpha_0}{\|\hat{\mu}^{(0)}\|} = \left\langle \frac{\hat{\mu}^{(0)}}{\|\hat{\mu}^{(0)}\|}, \eta \right\rangle \geq \frac{1}{\sqrt{d \log n_U}} \geq \frac{c'}{60\sqrt{d \log n_U}}$$

since $c' \leq 60$.

Now assume that $\alpha_t/\beta_t \geq c'/(60\sqrt{d \log n_U})$ and that $\alpha_0 \leq \alpha_t < c_5s$ for some $t \geq 0$. We aim to show that (b) holds for the same t , and start by controlling $e_1^\top f(\hat{\mu}^{(t)})$. Let $W = (W_1, \dots, W_d)^\top$ be an independent copy of Z_1 , independent of all other randomness in the problem, and define $u_t := \{\tanh(W_1\alpha_t + W_{-1}^\top \hat{\mu}_{-1}^{(t)}) - \tanh(W_1\alpha_t - W_{-1}^\top \hat{\mu}_{-1}^{(t)})\}/2$.

Then by applying the second part of Lemma 18 with $a = W_1\alpha_t$ and $b = W_{-1}^\top\hat{\mu}_{-1}^{(t)}$ (so that $a + b = W^\top\hat{\mu}^{(t)}$), we have

$$\begin{aligned} e_1^\top f(\hat{\mu}^{(t)}) &= \mathbb{E}\{W_1 \tanh(W_1\alpha_t + W_{-1}^\top\hat{\mu}_{-1}^{(t)})\} \\ &\geq \mathbb{E}\left\{\alpha_t W_1^2 - \frac{\alpha_t^3 W_1^4}{3} - \alpha_t W_1^2 (W_{-1}^\top\hat{\mu}_{-1}^{(t)})^2\right\} + \mathbb{E}(W_1 u_t) \\ &= \alpha_t(1 + s^2)(1 - \beta_t^2) - \alpha_t^3(1 + 2s^2 + s^4/3), \end{aligned}$$

where in the final step we have used the fact that u_t is an odd function of $W_{-1}^\top\hat{\mu}_{-1}^{(t)}$, which has a symmetric distribution about 0, conditional on $(W_1, \hat{\mu}_{-1}^{(t)})$, and hence $\mathbb{E}(W_1 u_t) = \mathbb{E}\{\mathbb{E}(W_1 u_t | W_1, \hat{\mu}_{-1}^{(t)})\} = 0$. From the assumption $s \geq c_4\zeta$, and using $\beta_t \leq 60(1+r)\zeta$ from Proposition 10, we have for sufficiently large c_4 that $\beta_t^2(1 + s^2) \leq \{60(1+r)\}^2(1 + r^2)s^2/c_4^2 \leq s^2/4$. By choosing $c_5 > 0$, depending only on r , sufficiently small, we may assume that $\alpha_t^2(1 + 2s^2 + s^4/3) < c_5^2(1 + 2r^2 + r^4/3)s^2 \leq s^2/4$. Recall the definition of f_{n_U} from (28). Since on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$, we have $\|\hat{\mu}^{(t)}\| \leq r + 3 \leq 2(r + \sqrt{d})$ as in the first line of the proof of Proposition 10, we have on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$ that

$$\begin{aligned} \alpha_{t+1} &= (1 - \gamma)e_1^\top f_{n_U}(\hat{\mu}^{(t)}) + \gamma e_1^\top \hat{\mu}_{n_L} \\ &\geq (1 - \gamma)\{e_1^\top f(\hat{\mu}^{(t)}) - \omega\|\hat{\mu}^{(t)}\|\} + \gamma(s - \phi) \\ &\geq (1 - \gamma)\{\alpha_t(1 + s^2/2) - \omega(\alpha_t + \beta_t)\} + \gamma(s - \phi). \end{aligned}$$

If $\gamma \geq \omega$, then $\phi \leq \omega\gamma^{-1/2} = \zeta \leq s/2$ (assuming $c_4 \geq 2$). If $\gamma < \omega$, then

$$\gamma\phi \leq \omega\gamma^{1/2} < \omega^{3/2} = \omega\zeta \leq \frac{\alpha_0\omega}{c'}\sqrt{d \log n_U} \leq \frac{2(1 - \gamma)\alpha_t\omega}{c'}\sqrt{d \log n_U}.$$

Hence, in either case, we have on $\Omega_1(\omega) \cap \Omega_2(\phi)$ that

$$\begin{aligned} \alpha_{t+1} &\geq (1 - \gamma)\alpha_t\{1 + s^2/2 - \omega(1 + 62/c')\sqrt{d \log n_U}\} + \frac{\gamma s}{2} \\ &\geq (1 - \gamma)\alpha_t\{1 + s^2/2 + (1 + 62/c')(2\gamma - \omega\sqrt{d \log n_U})\}, \end{aligned}$$

where the final bound holds provided we reduce c_5 to be at most $1/(2+124/c')$ if necessary.

Now, when $\gamma \leq \omega\sqrt{d \log n_U}$, we have $\zeta = \omega^{1/2} \wedge \omega\gamma^{-1/2} \geq \omega^{1/2}(d \log n_U)^{-1/4}$ and hence by the condition on s in the proposition, we have $s^2 \geq c_4^2\zeta^2 d \log n_U \geq c_4^2\omega\sqrt{d \log n_U}$. Thus, by increasing c_4 to be at least $\sqrt{4 + 248/c'}$ if necessary, we have $(1 + 62/c')\omega\sqrt{d \log n_U} \leq s^2/4$. Hence, in this case, and on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$,

$$\alpha_{t+1} \geq (1 - \gamma)(1 + s^2/4 + (2 + 124/c')\gamma)\alpha_t \geq \left(1 + \frac{s^2}{8}\right)\alpha_t \geq (1 + \omega\sqrt{d \log n_U})\alpha_t.$$

On the other hand, when $\gamma > \omega\sqrt{d \log n_U}$, we have

$$\alpha_{t+1} \geq (1 - \gamma)(1 + (1 + 62/c')\gamma)\alpha_t \geq (1 + \gamma)\alpha_t \geq (1 + \omega\sqrt{d \log n_U})\alpha_t.$$

Combining the two bounds above proves (b) for this given t .

It remains to verify (a) for $t + 1$, assuming that (a) and (b) hold up to and including t . Since $\beta_0 \leq \|\hat{\mu}^{(0)}\| \leq 60(\zeta \vee r\omega)$, we have by Proposition 10 that $\beta_{t+1} \leq 60(\zeta \vee r\omega) \leq 60(c')^{-1}\|\hat{\mu}^{(0)}\|$. Thus,

$$\frac{\alpha_{t+1}}{\beta_{t+1}} \geq \frac{\alpha_0}{60\|\hat{\mu}^{(0)}\|/c'} \geq \frac{c'}{60\sqrt{d \log n_U}},$$

which completes the induction. In particular, the geometric growth of α_t implied by (b) means that α_t will exceed $c_5 s$ for sufficiently large $t > 0$. This establishes our desired contradiction, and hence proves the result. \square

Proof of Proposition 5. Define $\phi_0 := \omega_0 \gamma^{-1/2}$, and recall the definitions of $\Omega_1(\omega)$ and $\Omega_2(\phi)$ from (30). By Proposition 9, there exists $C_r \geq 1$, depending only on r , such that for $\omega = C_r \omega_0$ and $\phi = C_r \phi_0$, we have $\mathbb{P}(\Omega_1(\omega) \cap \Omega_2(\phi)) \geq 1 - 2\delta$.

(i) By the definition of ω and ϕ , we have $\omega = \phi \gamma^{1/2}$. If we choose c such that $c \leq C_r^{-1} \min\{1/12, 1/(r+3)\}$, then $\omega \leq \min\{1/12, 1/(r+3)\}$. Thus, we may apply Propositions 10 and 11 to obtain that on $\Omega_1(\omega) \cap \Omega_2(\phi)$, we have

$$\limsup_{t \rightarrow \infty} \|\hat{\mu}^{(t)} - \mu^*\| \leq \limsup_{t \rightarrow \infty} (|\alpha_t| + \beta_t + \|\mu^*\|) \lesssim_r \zeta \vee \|\mu^*\|.$$

The first claim follows.

(ii) From Lemma 21 and by considering the case $d = 1$ separately, for the chosen η_0 , we have

$$\mathbb{P}(|e_1^\top \eta_0| \leq 1/\sqrt{d \log n_U}) \leq \sqrt{\frac{2}{\pi \log n_U}}.$$

Again, if we choose $c \leq C_r^{-1} \min\{1/12, 1/(r+3)\}$, then $\phi \gamma^{1/2} = \omega \leq \min\{1/12, 1/(r+3)\}$. Also, $\|\hat{\mu}^{(0)}\| = \zeta_0 \vee r\omega_0 \in [C_r^{-1}(\zeta \vee r\omega), \zeta \vee r\omega]$. Thus, applying Proposition 13 with $c' = 1/C_r$, there exists $c > 0$, depending only on r , and $t_0 \in \mathbb{N}$ such that on $\Omega_1(\omega) \cap \Omega_2(\phi) \cap \{|e_1^\top \eta_0| > 1/\sqrt{d \log n_U}\}$, we have $|\alpha_{t_0}| \geq cs$.

Since $\beta_0 \leq \|\hat{\mu}^{(0)}\| \leq \zeta \vee r\omega$, we can apply Proposition 10 to obtain that $\beta_t \leq 60(\zeta \vee r\omega)$ for all $t \geq 0$. Hence all conditions of Proposition 12 are satisfied, and the desired result then follows from (48). \square

To prove Theorem 6, we need the following proposition, which relates the loss of estimating μ^* to the operator norm loss of estimating $\mu^* \mu^{*\top}$.

Proposition 14. *Assume that $n \geq 3$ and that $(Z_1, Y_1, Y_1^*), \dots, (Z_n, Y_n, Y_n^*)$ are independent with*

$$Y_i^* \sim \text{Unif}\{-1, 1\}, Z_i | Y_i^* \sim \mathcal{N}_d(Y_i^* \mu^*, I_d), Y_i = Y_i^* \mathbb{1}_{\{i \leq n_L\}} \quad \text{for } i \in [n].$$

For $\mu \in \mathbb{R}^d$ and $i \in [n]$, let $L_i(\mu) := Y_i \mathbb{1}_{\{Y_i \neq 0\}} + \tanh\langle Z_i, \mu \rangle \mathbb{1}_{\{Y_i = 0\}}$, $\mu_{\text{tot}}(\mu) := \mu n^{-1} \sum_{i=1}^n L_i(\mu)$ and $\Sigma_b(\mu) := \mu \mu^\top - \mu_{\text{tot}}(\mu) \mu_{\text{tot}}(\mu)^\top$. For any $\delta \in (0, 1)$ and $B > 0$, we have with probability at least $1 - \delta$ that

$$\sup_{\mu: \|\mu\| \leq B} \{ \|\Sigma_b(\mu) - \mu^* \mu^{*\top}\|_{\text{op}} - (B+s)L(\mu, \mu^*) \} \lesssim B^2 (s^2 \vee 1) \left(\frac{d \log(2Bn + e) + \log(1/\delta)}{n} \right).$$

Proof. For any $\mu \in \mathbb{R}^d$ with $\|\mu\| \leq B$, we have

$$\begin{aligned} \|\Sigma_b(\mu) - \mu^* \mu^{*\top}\|_{\text{op}} &\leq \|\mu \mu^\top - \mu^* \mu^{*\top}\|_{\text{op}} + \|\mu_{\text{tot}}(\mu) \mu_{\text{tot}}(\mu)^\top\|_{\text{op}} \\ &\leq (B+s)L(\mu, \mu^*) + B^2 \left(\frac{1}{n} \sum_{i=1}^n L_i(\mu) \right)^2. \end{aligned}$$

Thus, it is enough to show that $\sup_{\mu: \|\mu\| \leq B} n^{-1} \sum_{i=1}^n L_i(\mu) \lesssim \sqrt{\frac{(s^2 \vee 1)\{d \log(2Bn+e) + \log(1/\delta)\}}{n}}$ with probability at least $1 - \delta$. To this end, we have

$$\sup_{\mu: \|\mu\| \leq B} \frac{1}{n} \sum_{i=1}^n L_i(\mu) = \frac{1}{n} \sum_{i=1}^{n_L} Y_i + \frac{1}{n} \sup_{\mu: \|\mu\| \leq B} \sum_{i=n_L+1}^n \tanh\langle Z_i, \mu \rangle. \quad (60)$$

For the first term on the right-hand side of (60), by Hoeffding's inequality, we have

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^{n_L} Y_i > \frac{\sqrt{2n_L \log(3/\delta)}}{n}\right) \leq \frac{\delta}{3}. \quad (61)$$

For the second term on the right-hand side of (60), let \mathcal{N} be a ε -net of $\{v : \|v\| \leq B\}$ with respect to the Euclidean distance, for some $\varepsilon \in (0, 1/2]$ to be specified later. Since a maximal ε -packing set is an ε -net, we may assume that $|\mathcal{N}| \leq (B + \varepsilon/2)^d / (\varepsilon/2)^d = (2B/\varepsilon + 1)^d$. Using the fact that $x \mapsto \tanh x$ is 1-Lipschitz, together with the Cauchy–Schwarz inequality, we have

$$\begin{aligned} & \sup_{v: \|v\| \leq B} \sum_{i=n_L+1}^n \tanh\langle Z_i, v \rangle \\ & \leq \sup_{v \in \mathcal{N}} \sum_{i=n_L+1}^n \tanh\langle Z_i, v \rangle + \sup_{u, v: \|u-v\| \leq \varepsilon} \sum_{i=n_L+1}^n (\tanh\langle Z_i, u \rangle - \tanh\langle Z_i, v \rangle) \\ & \leq \sup_{v \in \mathcal{N}} \sum_{i=n_L+1}^n \tanh\langle Z_i, v \rangle + \varepsilon \sum_{i=n_L+1}^n \|Z_i\|. \end{aligned}$$

Hence taking $\varepsilon = 1/n$, and defining $\tau := \frac{\log(3/\delta)}{d \log(2Bn+e)} > 0$, we have

$$\begin{aligned} & \mathbb{P}\left(\frac{1}{n} \sup_{v: \|v\| \leq B} \sum_{i=n_L+1}^n \tanh\langle Z_i, v \rangle > 2\sqrt{\frac{2(s^2 \vee 1)(1 + \tau)d \log(2Bn + e)}{n}}\right) \\ & \leq \mathbb{P}\left(\frac{1}{n} \sup_{v \in \mathcal{N}} \sum_{i=n_L+1}^n \tanh\langle Z_i, v \rangle > \sqrt{\frac{2(s^2 \vee 1)(1 + \tau)d \log(2Bn + e)}{n}}\right) \\ & \quad + \mathbb{P}\left(\frac{1}{n} \sum_{i=n_L+1}^n \|Z_i\| \geq \sqrt{2(s^2 \vee 1)(1 + \tau)nd \log(2Bn + e)}\right) \\ & \leq \frac{|\mathcal{N}|}{e^{(1+\tau)d \log(2Bn+e)}} + \mathbb{P}\left(\frac{1}{n} \sum_{i=n_L+1}^n \|Z_i\|^2 \geq 2(s^2 \vee 1)(1 + \tau)nd \log(2Bn + e)\right) \leq \frac{2\delta}{3}, \quad (62) \end{aligned}$$

where the penultimate bound uses Hoeffding's inequality and the Cauchy–Schwarz inequality and the final bound uses the fact that $\sum_{i=n_L+1}^n \|Z_i\|^2 \sim \chi_{n_U d}^2(n_U s^2)$ and [Birgé \(2001, Lemma 8.1\)](#). Combining (60), (61) and (62), we have with probability at least $1 - \delta$ that

$$\begin{aligned} \sup_{\mu: \|\mu\| \leq B} \frac{1}{n} \sum_{i=1}^n L_i(\mu) & \leq \frac{\sqrt{2n_L \log(3/\delta)}}{n} + 2\sqrt{\frac{2(s^2 \vee 1)(1 + \tau)d \log(2Bn + e)}{n}} \\ & \lesssim \sqrt{\frac{(s^2 \vee 1)\{d \log(2Bn + e) + \log(1/\delta)\}}{n}}, \end{aligned}$$

as desired. \square

Proof of Theorem 6. We write $\hat{\mu}_{[m]} \equiv \hat{\mu}_{[m]}^{(T)}$ for the T th (final) iterate of the EM update in Algorithm 3 starting from the m th random initializer $\hat{\mu}_{[m]}^{(0)}$. Let $\omega = C_r \omega_0$, $\phi = C_r \phi_0$, $\Omega_1(\omega)$ and $\Omega_2(\phi)$ be defined as in the proof of Proposition 5. Further, let $\Sigma_b(\mu)$ be defined as in Proposition 14. By Proposition 9, the first claim in the proof of Proposition 10 and Proposition 14, we have for some $C > 0$ depending only on r that

$$\mathbb{P} \left[\max_{m \in [M]} \sup_{T \in \mathbb{N}} \left\{ \|\Sigma_b(\hat{\mu}_{[m]}^{(T)}) - \mu^* \mu^{*\top}\|_{\text{op}} - (2r + 3)L(\hat{\mu}_{[m]}^{(T)}, \mu^*) \right\} \leq \frac{C\{d \log(rn) + \log(1/\delta)\}}{n} \right] \geq 1 - \delta.$$

In this balanced two-cluster setup, for the t th EM iteration starting from the m th random initializer, we have $-\hat{\mu}_1 = \hat{\mu}_2 = \hat{\mu}$, where we suppress the dependence on t and m for convenience. For $i \geq n_L + 1$, we have $L_{i,1} = e^{z_i^\top \hat{\mu}_1} / (e^{z_i^\top \hat{\mu}_1} + e^{z_i^\top \hat{\mu}_2})$, $L_{i,2} = e^{z_i^\top \hat{\mu}_2} / (e^{z_i^\top \hat{\mu}_1} + e^{z_i^\top \hat{\mu}_2})$ and hence $L_{i,2} - L_{i,1} = \tanh\langle Z_i, \hat{\mu} \rangle$. Thus, $\hat{\mu}_{\text{tot}} = n^{-1} \hat{\mu} \sum_{i=1}^n \{\tanh\langle Z_i, \hat{\mu} \rangle \mathbb{1}_{\{Y_i=0\}} + Y_i \mathbb{1}_{\{Y_i \neq 0\}}\}$. Also, we note that

$$\hat{\Sigma}_b = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^2 L_{i,k} (\hat{\mu}_k - \hat{\mu}_{\text{tot}}) (\hat{\mu}_k - \hat{\mu}_{\text{tot}})^\top = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^2 L_{i,k} \hat{\mu}_k \hat{\mu}_k^\top - \hat{\mu}_{\text{tot}} \hat{\mu}_{\text{tot}}^\top = \hat{\mu} \hat{\mu}^\top - \hat{\mu}_{\text{tot}} \hat{\mu}_{\text{tot}}^\top.$$

Consequently, using the notation of Proposition 14 and Algorithm 3, we have $\hat{Q} \equiv \hat{Q}^{(T)} \equiv \hat{Q}^{[\hat{m}]} = \Sigma_b(\hat{\mu}_{[\hat{m}]})$.

We consider two cases. If $\|\mu^*\| \leq \omega_0^{1/3} \wedge \zeta_0^{1/2}$, then by the proof of Proposition 5(i), we have on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$ that $\limsup_{T \rightarrow \infty} L(\hat{\mu}_{[m]}^{(T)}, \mu^*) \lesssim_r \zeta_0 \vee \|\mu^*\| \lesssim_r \omega_0^{1/3} \wedge \zeta_0^{1/2}$ for every $m \in [M]$. Thus, by Proposition 14, with probability at least $\mathbb{P}(\Omega_1(\omega) \cap \Omega_2(\phi)) - \delta \geq 1 - 3\delta$, we have

$$\begin{aligned} \limsup_{T \rightarrow \infty} \|\hat{Q} - \mu^* \mu^{*\top}\|_{\text{op}} &\lesssim_r (\omega_0^{1/3} \wedge \zeta_0^{1/2}) \limsup_{T \rightarrow \infty} \max_{m \in [M]} L(\hat{\mu}_{[m]}^{(T)}, \mu^*) \\ &\quad + (\omega_0^{2/3} \wedge \zeta_0) \left(\frac{d \log n + \log(1/\delta)}{n} \right) \\ &\lesssim_r \omega_0^{2/3} \wedge \zeta_0 = \frac{\omega_0}{\omega_0^{1/3} \wedge \zeta_0^{1/2}} \wedge \zeta_0 \leq \frac{\omega_0}{\|\mu^*\|} \wedge \zeta_0. \end{aligned} \quad (63)$$

We now turn to the case where $\|\mu^*\| > \omega_0^{1/3} \wedge \zeta_0^{1/2}$. Let \mathcal{M}_0 be the set of $m \in [M]$ such that $|\langle \hat{\mu}_{[m]}^{(0)}, \mu^* \rangle| / (\|\mu^*\| \|\hat{\mu}_{[m]}^{(0)}\|) \geq \sqrt{1/(d \log n_U)}$ and let $M_0 := |\mathcal{M}_0|$. By definition of the EM initializers, the random variables $\{\langle \hat{\mu}_{[m]}^{(0)}, \mu^* \rangle / (\|\mu^*\| \|\hat{\mu}_{[m]}^{(0)}\|) : m \in [M]\}$ are independent, and moreover, by Lemma 21 we have

$$\mathbb{P} \left(\frac{|\langle \hat{\mu}_{[m]}^{(0)}, \mu^* \rangle|}{\|\mu^*\| \|\hat{\mu}_{[m]}^{(0)}\|} \geq \sqrt{\frac{1}{d \log n_U}} \right) \geq 1 - \sqrt{\frac{2}{\pi \log n_U}} > \frac{3}{5}.$$

Defining $\Omega_3 := \{M_0 > M/2\}$, by Hoeffding's inequality, we have

$$\mathbb{P}(\Omega_3^c) \leq e^{-M/50}.$$

Let

$$\mathcal{M}_1 := \{m \in [M] \setminus \{\hat{m}\} : \|\hat{Q}^{[m]} - \hat{Q}^{[\hat{m}]}\|_{\text{op}} \leq \text{median}(\|\hat{Q}^{[m']} - \hat{Q}^{[\hat{m}]}\|_{\text{op}} : m' \in [M] \setminus \{\hat{m}\})\}.$$

Since $|\mathcal{M}_1 \cup \{\hat{m}\}| \geq \lceil (M-1)/2 \rceil + 1 > M/2$, we have on Ω_3 that $\mathcal{M}_0 \cap (\mathcal{M}_1 \cup \{\hat{m}\}) \neq \emptyset$. Thus, on the event $\Omega_1(\omega) \cap \Omega_2(\phi) \cap \Omega_3$, we can let $\tilde{m} := \min(\mathcal{M}_0 \cap (\mathcal{M}_1 \cup \{\hat{m}\}))$, so by definition of \hat{m} , we have

$$\begin{aligned} \|\hat{Q}^{[\hat{m}]} - \mu^* \mu^{*\top}\|_{\text{op}} &\leq \|\hat{Q}^{[\hat{m}]} - \hat{Q}^{[\tilde{m}]}\|_{\text{op}} + \|\hat{Q}^{[\tilde{m}]} - \mu^* \mu^{*\top}\|_{\text{op}} \\ &\leq \text{median}(\|\hat{Q}^{[m']} - \hat{Q}^{[\tilde{m}]}\|_{\text{op}} : m' \in [M] \setminus \{\hat{m}\}) + \|\hat{Q}^{[\tilde{m}]} - \mu^* \mu^{*\top}\|_{\text{op}} \\ &\leq \text{median}(\|\hat{Q}^{[m']} - \hat{Q}^{[\tilde{m}]}\|_{\text{op}} : m' \in [M] \setminus \{\tilde{m}\}) + \|\hat{Q}^{[\tilde{m}]} - \mu^* \mu^{*\top}\|_{\text{op}} \\ &\leq \max_{m, m' \in \mathcal{M}_0} \|\hat{Q}^{[m]} - \hat{Q}^{[m']}\|_{\text{op}} + \|\hat{Q}^{[\tilde{m}]} - \mu^* \mu^{*\top}\|_{\text{op}} \\ &\leq 3 \max_{m \in \mathcal{M}_0} \|\hat{Q}^{[m]} - \mu^* \mu^{*\top}\|_{\text{op}}. \end{aligned}$$

Since $\omega_0 \leq (d \log n)^{-3}$, by discussing cases of $\gamma < \omega$, $\omega \leq \gamma \leq \omega^{2/3}$ and $\gamma > \omega^{2/3}$, we see that $\omega_0^{1/3} \wedge \zeta_0^{1/2} \geq \zeta_0 \sqrt{d \log n}$. From the proof of Proposition 5(ii), we have on the event $\Omega_1(\omega) \cap \Omega_2(\phi)$ that $\limsup_{T \rightarrow \infty} \max_{m \in \mathcal{M}_0} L(\hat{\mu}_{[m]}^{(T)}, \mu^*) \lesssim_r \frac{\omega_0}{\|\mu^*\|} \wedge (\omega_0 \gamma^{-1/2})$. Let Ω_4 be the event on which the conclusion of Proposition 14 holds. Then on $\Omega_1(\omega) \cap \Omega_2(\phi) \cap \Omega_3 \cap \Omega_4$, we therefore have

$$\begin{aligned} \limsup_{T \rightarrow \infty} \|\hat{Q}^{[\hat{m}]} - \mu^* \mu^{*\top}\|_{\text{op}} &\lesssim_r \limsup_{T \rightarrow \infty} \max_{m \in \mathcal{M}_0} L(\hat{\mu}_{[m]}^{(T)}, \mu^*) + \frac{d \log n + \log(1/\delta)}{n} \\ &\lesssim_r \left(\frac{\omega_0}{\|\mu^*\|} \wedge \frac{\omega_0}{\gamma^{1/2}} \right) + \frac{d \log n + \log(1/\delta)}{n} \\ &\lesssim_r \frac{\omega_0}{\|\mu^*\|} \wedge \zeta_0. \end{aligned} \tag{64}$$

The desired result follows by combining (63) and (64), and the fact that $\mathbb{P}(\Omega_1(\omega) \cap \Omega_2(\phi) \cap \Omega_3 \cap \Omega_4) \geq 1 - 3\delta - e^{-M/50}$. \square

5.5 Proof of Corollary 7

Proof of Corollary 7. Fix $P \in \mathcal{P}_d$, define $Z_i := PX_i$ for $i \in [n]$, $\mu^* := P\nu^*$, $\delta := \varepsilon/\{4\binom{p}{d}\}$, $\omega_0 := \sqrt{\frac{d \log n + \log(1/\delta)}{n_U}}$ and $\zeta_0 := \omega_0^{1/2} \wedge \omega_0 \gamma^{-1/2}$. Then, provided $C_1 > 2$, we have $\delta \geq 2e^{-n/2}/p^d > 2e^{-n}$, and $\|\mu^*\| \leq \|\nu^*\| \leq r$. Let $c > 0$ be chosen, depending only on r , to satisfy Theorem 6. By increasing $C_1 > 0$, depending only on r , if necessary, we may assume that $\omega_0 \leq \min\{c, (d \log n)^{-3}\}$. Hence, since $(P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top = \mu^* \mu^{*\top}$, we can apply Theorem 6 to obtain that for some $C'_2 > 0$ depending only on r , with probability at least $1 - 3\delta - e^{-M/50}$ we have that

$$\begin{aligned} \limsup_{T \rightarrow \infty} \|\psi^{(T)}((PX_i, Y_i)_{i \in [n]}) - (P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top\|_{\text{op}} &\leq C'_2 \zeta_0 \\ &\leq C_2 \min \left[\left\{ \frac{d \log(p \vee n) + \log(1/\varepsilon)}{n} \right\}^{1/4}, \sqrt{\frac{d \log(p \vee n) + \log(1/\varepsilon)}{n_L}} \right] \leq \frac{(\nu_{\min}^*)^2}{4}. \end{aligned}$$

Since $\psi^{(T)}$ is permutation equivariant for each $T \geq 0$, by Fatou's lemma and a union bound, we have that

$$\begin{aligned} & \limsup_{T \rightarrow \infty} \mathbb{P} \left(\max_{P \in \mathcal{P}_d} \left\| \psi^{(T)} \left((PX_i, Y_i)_{i \in [n]} \right) - (P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top \right\|_{\text{op}} > \frac{(\nu_{\min}^*)^2}{4} \right) \\ & \leq \mathbb{P} \left(\limsup_{T \rightarrow \infty} \max_{P \in \mathcal{P}_d} \left\| \psi^{(T)} \left((PX_i, Y_i)_{i \in [n]} \right) - (P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top \right\|_{\text{op}} > \frac{(\nu_{\min}^*)^2}{4} \right) \\ & \leq \sum_{P \in \mathcal{P}_d} \mathbb{P} \left(\limsup_{T \rightarrow \infty} \left\| \psi^{(T)} \left((PX_i, Y_i)_{i \in [n]} \right) - (P\Sigma_w P^\top)^{-1} P\Sigma_b P^\top \right\|_{\text{op}} > \frac{(\nu_{\min}^*)^2}{4} \right) \\ & \leq \binom{p}{d} (3\delta + e^{-M/50}) \leq \frac{3}{4}\varepsilon + e^{-M/50+d \log p} \leq \varepsilon. \end{aligned}$$

The result now follows from Theorem 2, noting that $\gamma_{\min} = (\nu_{\min}^*)^2$ and $\gamma_{\max} = (\nu_{\max}^*)^2$. \square

6 Auxiliary lemmas

Lemma 15. *Suppose that $K = 2$ and \mathcal{C} is defined as in (11). Let $(-\hat{\mu}^{(t)}, \hat{\mu}^{(t)}, I_d) \in \mathcal{C}$ be the t th iterate of the EM iteration described in (4) and (5) with data $(Z_1, Y_1), \dots, (Z_n, Y_n)$, starting from $(-\hat{\mu}^{(0)}, \hat{\mu}^{(0)}, I_d)$. Then for all $t \geq 1$, we have*

$$\hat{\mu}^{(t)} = \frac{1}{n} \left\{ \sum_{i: Y_i \neq 0} (-1)^{Y_i} Z_i + \sum_{i: Y_i = 0} Z_i \tanh \langle Z_i, \hat{\mu}^{(t-1)} \rangle \right\}.$$

Proof. At step $t \geq 1$, in the E-step, by (4), we have for $k \in \{1, 2\}$ that $L_{i,k} = \mathbb{1}_{\{Y_i=k\}}$ if $Y_i \neq 0$ and

$$L_{i,k} = \frac{e^{-\|Z_i - (-1)^k \hat{\mu}^{(t-1)}\|^2/2}}{e^{-\|Z_i - \hat{\mu}^{(t-1)}\|^2/2} + e^{-\|Z_i + \hat{\mu}^{(t-1)}\|^2/2}}$$

otherwise. In the M-step, defining

$$Q(\mu \mid \hat{\mu}^{(t-1)}) := \frac{1}{n} \sum_{i=1}^n (L_{i,1} \|Z_i + \mu\|^2 + L_{i,2} \|Z_i - \mu\|^2),$$

we have $\hat{\mu}^{(t)} = \operatorname{argmin}_{\mu \in \mathbb{R}^d} Q(\mu \mid \hat{\mu}^{(t-1)})$. Differentiating $Q(\mu \mid \hat{\mu}^{(t-1)})$ with respect to μ , we obtain

$$\hat{\mu}^{(t)} = \frac{1}{n} \sum_{i=1}^n (L_{i,2} - L_{i,1}) Z_i.$$

The desired result follows since $L_{i,2} - L_{i,1} = (-1)^{Y_i}$ if $Y_i \in \{1, 2\}$, and

$$L_{i,2} - L_{i,1} = \frac{e^{-\|Z_i - \hat{\mu}^{(t-1)}\|^2/2} - e^{-\|Z_i + \hat{\mu}^{(t-1)}\|^2/2}}{e^{-\|Z_i - \hat{\mu}^{(t-1)}\|^2/2} + e^{-\|Z_i + \hat{\mu}^{(t-1)}\|^2/2}} = \frac{e^{\langle Z_i, \hat{\mu}^{(t-1)} \rangle} - e^{-\langle Z_i, \hat{\mu}^{(t-1)} \rangle}}{e^{\langle Z_i, \hat{\mu}^{(t-1)} \rangle} + e^{-\langle Z_i, \hat{\mu}^{(t-1)} \rangle}} = \tanh \langle Z_i, \hat{\mu}^{(t-1)} \rangle$$

if $Y_i = 0$. \square

Lemma 16. *Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P$ for some distribution P on \mathbb{R}^d . If $\|\mu^*\| \leq n^{-1/4}$, then for any Borel measurable function $\psi : (\mathbb{R}^d)^n \rightarrow \{0, 1\}$ of the null hypothesis $H_0 : P = \mathcal{N}_d(0, I_d)$ against the alternative $H_1 : P = \frac{1}{2}\mathcal{N}_d(\mu^*, I_d) + \frac{1}{2}\mathcal{N}_d(-\mu^*, I_d)$, we have*

$$\mathbb{P}_{H_0}(\psi(X_1, \dots, X_n) = 1) + \mathbb{P}_{H_1}(\psi(X_1, \dots, X_n) = 0) > 1/2.$$

Proof. Write $X = (X_1, \dots, X_n)^\top$. Observe that, writing d_{TV} for the total variation distance between probability measures,

$$\begin{aligned} \mathbb{P}_{H_0}(\psi(X) = 1) + \mathbb{P}_{H_1}(\psi(X) = 0) &\geq 1 - d_{\text{TV}}(\mathbb{P}_{H_0}, \mathbb{P}_{H_1}) \\ &= 1 - \frac{1}{2} \int \left| \frac{d\mathbb{P}_{H_1}}{d\mathbb{P}_{H_0}} - 1 \right| d\mathbb{P}_{H_0} \geq 1 - \frac{1}{2} \left\{ \int \left(\frac{d\mathbb{P}_{H_1}}{d\mathbb{P}_{H_0}} - 1 \right)^2 d\mathbb{P}_{H_0} \right\}^{1/2} \\ &= 1 - \frac{1}{2} \left\{ \int \left(\frac{d\mathbb{P}_{H_1}}{d\mathbb{P}_{H_0}} \right)^2 d\mathbb{P}_{H_0} - 1 \right\}^{1/2}. \end{aligned} \quad (65)$$

To control the chi-squared divergence in the right-hand side of (65) above, we let $\xi = (\xi_1, \dots, \xi_n)^\top$ have independent Rademacher components and $W = (W_{i,j})_{i \in [n], j \in [d]}$ be a random matrix with independent $N(0, 1)$ entries, independent of ξ . Then $X \stackrel{d}{=} W$ under H_0 and $X \stackrel{d}{=} \xi \mu^{*\top} + W$ under H_1 . Let $\tilde{\xi}$ be an independent copy of ξ . Using the Ingster–Suslina device, see, e.g., Ingster and Suslina (2012), Liu, Gao and Samworth (2021), Lemma 21, we have that

$$\int \left(\frac{d\mathbb{P}_{H_1}}{d\mathbb{P}_{H_0}} \right)^2 d\mathbb{P}_{H_0} = \mathbb{E} \exp \langle \xi \mu^{*\top}, \tilde{\xi} \mu^{*\top} \rangle = \cosh^n(\|\mu^*\|^2) \leq e^{n\|\mu^*\|^4/2} \leq e^{1/2},$$

where we used the fact that $\cosh x \leq e^{x^2/2}$ for all $x \in \mathbb{R}$ in the penultimate step. The desired result follows from substituting the above bound into (65) and the fact that $1 - (e^{1/2} - 1)^{1/2}/2 > 1/2$. \square

We prove a generalization of Cochran’s theorem for quadratic forms of independent Gaussian random vectors with a common covariance matrix, which result in independent noncentral Wishart distributions. Recall that if X is a matrix, then $\text{vec}(X)$ is the vectorization of X , obtained by stacking its columns on top of each other. The Kronecker product between matrices $A = (A_{i,j})_{i \in [m], j \in [n]}$ and B is defined as

$$A \otimes B := \begin{pmatrix} A_{1,1}B & \cdots & A_{1,n}B \\ \vdots & \ddots & \vdots \\ A_{m,1}B & \cdots & A_{m,n}B \end{pmatrix}.$$

Recall also that when $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}_d(0, \Sigma)$, the matrix $\sum_{i=1}^n X_i X_i^\top$ has a d -dimensional Wishart distribution with n degrees of freedom and covariance matrix $\Sigma \in \mathbb{S}^{d \times d}$, denoted $\mathcal{W}_d(n, \Sigma)$. More generally, $\sum_{i=1}^n (X_i + \mu_i)(X_i + \mu_i)^\top$ has a non-central Wishart distribution with n degrees of freedom, covariance matrix Σ and non-centrality matrix $\Omega = \sum_{i=1}^n \mu_i \mu_i^\top$, written $\mathcal{W}_d(n, \Sigma; \Omega)$. Thus $\mathcal{W}_d(n, \Sigma; 0) \stackrel{d}{=} \mathcal{W}_d(n, \Sigma)$.

Lemma 17. *Let Z_1, \dots, Z_n be independent with $Z_i \sim \mathcal{N}_d(\mu_i, \Sigma)$ for $i \in [n]$, and write $Z := (Z_1, \dots, Z_n)^\top \in \mathbb{R}^{n \times d}$ and $M := \mathbb{E}(Z)$. If $P_1, \dots, P_k \in \mathbb{R}^{n \times n}$ are positive semidefinite matrices such that $P_1 + \dots + P_k = I_n$ and $\text{rank}(P_1) + \dots + \text{rank}(P_k) = n$, then $Z^\top P_1 Z, \dots, Z^\top P_k Z$ are independent with $Z^\top P_r Z \sim \mathcal{W}_d(\text{rank}(P_r), \Sigma; M^\top P_r M)$.*

Proof. As in the proof of the classical Cochran’s theorem (Cochran, 1934), we first note that P_1, \dots, P_k can be simultaneously diagonalised such that

$$P_r = Q D_r Q^\top,$$

for some $Q \in \mathbb{O}^{n \times n}$ and $D_r = \text{diag}((\mathbb{1}_{\{j \in S_r\}})_{j \in [n]})$, where $S_r \subseteq [n]$, $|S_r| = \text{rank}(P_r)$ and $S_r \cap S_{r'} = \emptyset$ for all $r \neq r'$. In particular, P_1, \dots, P_k satisfy $P_r^2 = P_r$ for $r \in [k]$ and $P_r P_{r'} = 0$ for all $r \neq r'$. Since $P_1 Z, \dots, P_k Z$ are jointly Gaussian, with

$$\begin{aligned} \text{Cov}(\text{vec}(P_r Z), \text{vec}(P_{r'} Z)) &= \text{Cov}((I_d \otimes P_r) \text{vec}(Z), (I_d \otimes P_{r'}) \text{vec}(Z)) \\ &= (I_d \otimes P_r)(\Sigma \otimes I_n)(I_d \otimes P_{r'})^\top = 0, \end{aligned}$$

we have that $P_1 Z, \dots, P_k Z$ are independent. But $Z^\top P_r Z = (P_r Z)^\top P_r Z$, so it follows that $Z^\top P_1 Z, \dots, Z^\top P_k Z$ are independent. Moreover, writing $W = (W_1, \dots, W_n)^\top := Q^\top Z$, we have $\text{vec}(Z) \sim \mathcal{N}_{nd}(\text{vec}(M), \Sigma \otimes I_n)$, so

$$\text{vec}(W) = (I_d \otimes Q^\top) \text{vec}(Z) \sim \mathcal{N}_{nd}((I_d \otimes Q^\top) \text{vec}(M), \Sigma \otimes I_n) \stackrel{d}{=} \mathcal{N}_{nd}(\text{vec}(Q^\top M), \Sigma \otimes I_n).$$

Therefore,

$$\begin{aligned} Z^\top P_r Z &= W^\top D_r W = \sum_{i \in S_r} W_i W_i^\top \sim \mathcal{W}_d(|S_r|, \Sigma; \sum_{i \in S_r} \mathbb{E}(W_i) \mathbb{E}(W_i)^\top) \\ &\stackrel{d}{=} \mathcal{W}_d(\text{rank}(P_r), \Sigma; M^\top P_r M), \end{aligned}$$

as desired. □

Lemma 18. *For any $a, b \in \mathbb{R}$, we have*

$$\frac{1}{2} \{ \tanh(a+b) - \tanh(a-b) \} \leq |b|$$

and

$$\frac{a}{2} \{ \tanh(a+b) + \tanh(a-b) \} \geq a^2 - \frac{a^4}{3} - a^2 b^2$$

Proof. For the first inequality, since the left-hand side is an increasing function of b , and an even function of a , we may assume that $a \geq 0$ and $b \geq 0$. Notice that $\frac{\partial}{\partial a} (\tanh(a+b) - \tanh(a-b)) = 1/\cosh^2(a+b) - 1/\cosh^2(|a-b|) \leq 0$, since $x \mapsto \cosh(x)$ is an increasing function on $[0, \infty)$. Hence

$$\frac{1}{2} \{ \tanh(a+b) - \tanh(a-b) \} \leq \tanh b \leq b.$$

as desired.

For the second inequality, since both sides are even functions of both a and b , we may again assume without loss of generality that $a > 0$ and $b \geq 0$. We may also assume that $b \leq 1$ since otherwise, the right-hand side is negative and the inequality holds trivially. But then

$$\begin{aligned} \frac{1}{2} \{ \tanh(a+b) + \tanh(a-b) \} &= \frac{1}{2} \left(\frac{\tanh a + \tanh b}{1 + \tanh a \tanh b} + \frac{\tanh a - \tanh b}{1 - \tanh a \tanh b} \right) \\ &= \frac{\tanh a}{(1 - \tanh^2 a \tanh^2 b) \cosh^2 b} \geq \frac{\tanh a}{\cosh^2 b} \\ &\geq (1 - b^2) \tanh a \geq (1 - b^2) \left(a - \frac{a^3}{3} \right) \geq a - \frac{a^3}{3} - ab^2, \end{aligned}$$

as desired. Here, the second inequality holds because $(1 - b^2) \cosh^2 b \leq (1 - b^2) e^{b^2} \leq 1$. □

Lemma 19. Let $H : [0, \infty) \rightarrow [0, \infty)$ be an increasing, concave function with $H'(x_0) < 1$ for some $x_0 \geq 0$ and either $H(0) > 0$ or both $H(0) = 0$ and $H'(0) > 1$. Then there exists a unique $\alpha^* > 0$ such that

$$H(\alpha) - \alpha \begin{cases} > 0 & \alpha \in (0, \alpha^*) \\ = 0 & \alpha = \alpha^* \\ < 0 & \alpha \in (\alpha^*, \infty). \end{cases}$$

Moreover, if $\alpha_0 > 0$, then the sequence $(\alpha_t)_{t \geq 0}$ given by $\alpha_t := H(\alpha_{t-1})$ monotonically converges to α^* .

Proof. For the first claim, consider the concave function $\tilde{H}(x) := H(x) - x$, which satisfies $\tilde{H}(x) > 0$ for sufficiently small $x > 0$, and for $x \geq x_0$, we have that any supergradient $v_x \in \mathbb{R}$ of \tilde{H} at x satisfies $v_x \leq -\{1 - H'(x_0)\} < 0$. It follows that $\tilde{H}(x) \rightarrow -\infty$ as $x \rightarrow \infty$, so by the intermediate value theorem, there exists $\alpha^* \in (0, \infty)$ such that $\tilde{H}(\alpha^*) = 0$, i.e. $H(\alpha^*) = \alpha^*$. Again using the facts that $\tilde{H}(x) > 0$ for sufficiently small $x > 0$, and $\tilde{H}(x) \rightarrow -\infty$ as $x \rightarrow \infty$, we see that the concave function \tilde{H} can only cross the x -axis at one positive value α^* , and $\tilde{H}(\alpha) > 0$ for $\alpha \in (0, \alpha^*)$ and $\tilde{H}(\alpha) < 0$ for $\alpha \in (\alpha^*, \infty)$.

Next, note that if $\alpha \in (0, \alpha^*)$, then $\alpha < H(\alpha) < H(\alpha^*) = \alpha^*$. Thus, if $\alpha_0 < \alpha^*$, then $(\alpha_t)_{t \geq 0}$ is an increasing sequence, bounded above by α^* , so it converges to a limit. But then, taking limits on both sides of the recursion $\alpha_t := H(\alpha_{t-1})$, we deduce that this limit must be α^* . A similar argument can be used to show that if $\alpha_0 \in (\alpha^*, \infty)$ then $(\alpha_t)_{t \geq 0}$ decreases down to the limit α^* , while if $\alpha_0 = \alpha^*$, then $\alpha_t = \alpha^*$ for all t . \square

Lemma 20. Let μ^* be a non-zero vector in \mathbb{R}^d , let $Z \sim \frac{1}{2}\mathcal{N}_d(-\mu^*, I_d) + \frac{1}{2}\mathcal{N}_d(\mu^*, I_d)$, let $\eta := \mu^*/\|\mu^*\|$, and define $q : [0, \infty) \rightarrow [0, \infty)$ by

$$q(\alpha) := \begin{cases} \alpha^{-1} \eta^\top \mathbb{E}(Z \tanh\langle \alpha \eta, Z \rangle) & \text{if } \alpha > 0 \\ 1 + \|\mu^*\|^2 & \text{if } \alpha = 0. \end{cases}$$

Then q is a differentiable function with $q(\|\mu^*\|) = 1$ and for any $h \geq \|\mu^*\|$, we have

$$\sup_{\alpha \in [0, h]} \frac{q'(\alpha)}{\alpha} \leq -\frac{e^{-h^2/2}}{3 \cdot 2^{11} \sqrt{2\pi} (h^5 \vee 1)}.$$

Proof. Write $s := \|\mu^*\|$. The fact that $q(s) = 1$ follows from [Xu, Hsu and Maleki \(2016, Theorem 1\)](#). By [Wu and Zhou \(2022, Lemma 3\(4\)\)](#), q is differentiable with $q'(\alpha) \leq -(2\alpha/3) \cdot \mathbb{E}(\tilde{Z}^4 / \cosh^2(\alpha \tilde{Z}))$ for $\alpha \in [0, \infty)$, where $\tilde{Z} \sim \frac{1}{2}\mathcal{N}(-s, 1) + \frac{1}{2}\mathcal{N}(s, 1)$. We can

now compute that for $\alpha, s \in [0, h]$,

$$\begin{aligned}
\mathbb{E}\left(\frac{\tilde{Z}^4}{\cosh^2(\alpha\tilde{Z})}\right) &\geq \mathbb{E}(\tilde{Z}^4 e^{-2\alpha|\tilde{Z}|}) \geq \frac{1}{2\sqrt{2\pi}} \int_0^\infty y^4 e^{-2\alpha y} e^{-(y-s)^2/2} dy \\
&= \frac{1}{2\sqrt{2\pi}} \int_0^\infty y^4 e^{-(y-s+2\alpha)^2/2-2\alpha s+2\alpha^2} dy \\
&\geq \begin{cases} \frac{e^{-2\alpha s+2\alpha^2}}{2\sqrt{2\pi}} \int_0^{\frac{1}{2(2\alpha-s)}} \frac{1}{2} y^4 e^{-(2\alpha-s)^2/2} dy & \text{if } s+1 \leq 2\alpha \\ \frac{e^{-2\alpha s+2\alpha^2}}{2\sqrt{2\pi}} \int_{s-2\alpha+1}^{s-2\alpha+2} y^4 e^{-2} dy & \text{if } s+1 > 2\alpha \end{cases} \\
&\geq \begin{cases} \frac{e^{-s^2/2}}{2^7 \cdot 5\sqrt{2\pi}} \left(\frac{1}{2\alpha-s}\right)^5 & \text{if } s+1 \leq 2\alpha \\ \frac{e^{-2\alpha s+2\alpha^2-2}}{10\sqrt{2\pi}} & \text{if } s+1 > 2\alpha \end{cases} \\
&\geq \frac{e^{-h^2/2}}{2^{12} \cdot 5\sqrt{2\pi}(h^5 \vee 1)},
\end{aligned}$$

which establishes the desired bound. \square

Lemma 21. *Let $d \geq 2$, and let $\eta = (\eta_1, \dots, \eta_d)^\top \sim \text{Unif}(\mathbb{S}^{d-1})$. Then for any $a > 0$, we have*

$$\mathbb{P}\left(|\eta_1| \leq \frac{a}{\sqrt{d}}\right) \leq \sqrt{\frac{2}{\pi}} a.$$

Proof. Letting $Z = (Z_1, \dots, Z_d)^\top \sim N_d(0, I_d)$, we have $\eta \stackrel{d}{=} Z/\|Z\|$ and in particular $\eta_1^2 \stackrel{d}{=} Z_1^2/(Z_1^2 + \dots + Z_d^2) \sim \text{Beta}(1/2, (d-1)/2)$. Thus,

$$\begin{aligned}
\mathbb{P}\left(|\eta_1| \leq \frac{a}{\sqrt{d}}\right) &= \mathbb{P}\left(\frac{Z_1^2}{\|Z\|^2} \leq \frac{a^2}{d}\right) = \frac{\Gamma(d/2)}{\Gamma(1/2)\Gamma((d-1)/2)} \int_0^{a^2/d} t^{-1/2}(1-t)^{(d-3)/2} dt \\
&\leq \frac{2\Gamma(d/2)a}{\sqrt{d}\Gamma(1/2)\Gamma((d-1)/2)} \leq \sqrt{\frac{2}{\pi}} a,
\end{aligned}$$

where the final bound uses, e.g., [Dümbgen, Samworth and Wellner \(2021, Corollary 11\)](#). \square

References

- Ahfock, D. C., Astle, W. J. and Richardson, S. (2021) Statistical properties of sketching algorithms. *Biometrika*, **108**, 283–297.
- Akçay, S., Atapour-Abarghouei, A. and Breckon, T. P. (2019) Ganomaly: Semi-supervised anomaly detection via adversarial training. In *14th Asian Conference on Computer Vision, Revised Selected Papers, Part III 14*, 622–637, Springer.
- Alon, U., Barkai, N., Notterman, D. A., Gish, K., Ybarra, S., Mack, D. and Levine, A. J. (1999) Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc. Nat. Acad. Sci.*, **96**, 6745–6750.

- Anderlucci, L., Fortunato, F. and Montanari, A. (2022) High-dimensional clustering via random projections. *J. Classification*, **39**, 191–216.
- Anderson, T. W. (2003) *An Introduction to Multivariate Statistical Analysis*. Wiley Series in Probability and Statistics.
- Azizyan, M., Singh, A. and Wasserman, L. (2013) Minimax theory for high-dimensional Gaussian mixtures with sparse mean separation. In *Adv. Neur. Inf. Proc. Syst.*, 2139–2147.
- Azizyan, M., Singh, A. and Wasserman, L. (2015) Efficient sparse clustering of high-dimensional non-spherical Gaussian mixtures. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics*, 37–45.
- Balakrishnan, S., Wainwright, M. J. and Yu, B. (2017) Statistical guarantees for the EM algorithm: from population to sample-based analysis. *Ann. Statist.*, **45**, 77–120.
- Bingham, E. and Mannila, H. (2001) Random projection in dimensionality reduction: applications to image and text data. In *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 245–250.
- Birgé, L. (2001) An alternative point of view on Lepski’s method. In de Gunst, M., Klaassen, C. and van der Vaart, A. Eds., *Lecture Notes-Monograph Series 36*, 113–133.
- Boucheron, S., Lugosi, G. and Massart, P. (2013) *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press.
- Butler, A., Hoffman, P., Smibert, P., Papalexi, E. and Satija, R. (2018) Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotech.*, **36**, 411–420.
- Cai, T. T. and Liu, W. (2011) A direct estimation approach to sparse linear discriminant analysis. *J. Amer. Statist. Assoc.*, **106**, 1566–1577.
- Cai, T. T. and Zhang, L. (2019) High dimensional linear discriminant analysis: optimality, adaptive algorithm and missing data. *J. Roy. Statist. Soc., Ser. B*, **81**, 675–705.
- Cannings, T. I. (2021) Random projections: Data perturbation for classification problems. *Wiley Interdisciplinary Reviews: Computational Statistics*, **13**, e1499.
- Cannings, T. I., Berrett, T. B. and Samworth, R. J. (2020) Local nearest neighbour classification with applications to semi-supervised learning. *Ann. Statist.*, **48**, 1789–1814.
- Cannings, T. I. and Samworth, R. J. (2017) Random-projection ensemble classification. *J. Roy. Statist. Soc., Ser. B*, **79**, 959–1035.
- Chakraborty, A. and Cai, T. (2018) Efficient and adaptive linear regression in semi-supervised settings. *Ann. Statist.*, **46**, 1541–1572.
- Chapelle, O., Schölkopf, B. and Zien, A. (Eds.) (2006) *Semi-Supervised Learning*. The MIT Press.

- Cheplygina, V., de Bruijne, M. and Pluim, J. P. (2019) Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis. *Medical Image Analysis*, **54**, 280–296.
- Cochran, W. G. (1934) The distribution of quadratic forms in a normal system, with applications to the analysis of covariance. *Mathematical Proceedings of the Cambridge Philosophical Society*, **30**, 178–191.
- Dasgupta, S. (1999) Learning mixtures of Gaussians. In *The 40th Annual Symposium on Foundations of Computer Science*, 634–644.
- Dasgupta, S. and Gupta, A. (2003) An elementary proof of a theorem of Johnson and Lindenstrauss. *Random Structures & Algorithms*, **22**, 60–65.
- Daskalakis, C., Tzamos, C. and Zampetakis, M. (2017) Ten steps of EM suffice for mixtures of two Gaussians. In *Conference on Learning Theory*, 704–710, PMLR.
- Davis, D., Diaz, M. and Wang, K. (2021) Clustering a mixture of Gaussians with unknown covariance. *arXiv preprint arXiv:2110.01602*.
- de Souto, M. C., Costa, I. G., de Araujo, D. S., Ludermir, T. B. and Schliep, A. (2008) Clustering cancer gene expression data: a comparative study. *BMC Bioinformatics*, **9**, 497.
- Devroye, L., Györfi, L. and Lugosi, G. (2013) *A Probabilistic Theory of Pattern Recognition*, vol. 31. Springer Science & Business Media.
- Dimitriadou, E., Weingessel, A. and Hornik, K. (2002) A combination scheme for fuzzy clustering. *International Journal of Pattern Recognition and Artificial Intelligence*, **16**, 901–912.
- Dobriban, E. and Liu, S. (2019) Asymptotics for sketching in least squares regression. In *Adv. Neur. Inf. Proc. Syst.*, 3675–3685.
- Doss, N., Wu, Y., Yang, P. and Zhou, H. H. (2023) Optimal estimation of high-dimensional Gaussian mixtures. *Ann. Statist.*, **51**, 62–95.
- Dümbgen, L., Samworth, R. J. and Wellner, J. A. (2021) Bounding distributional errors via density ratios. *Bernoulli*, **27**, 818–852.
- Durrant, R. J. and Kabán, A. (2015) Random projections as regularizers: learning a linear discriminant from fewer observations than dimensions. *Machine Learning*, **99**, 257–286.
- Dwivedi, R., Ho, N., Khamaru, K., Wainwright, M., Jordan, M. and Yu, B. (2020a) Sharp analysis of expectation-maximization for weakly identifiable models. In *International Conference on Artificial Intelligence and Statistics*, 1866–1876, PMLR.
- Dwivedi, R., Ho, N., Khamaru, K., Wainwright, M. J., Jordan, M. I. and Yu, B. (2020b) Singularity, misspecification and the convergence rate of EM. *Ann. Statist.*, **48**, 3161–3182.

- Eisen, M. B., Spellman, P. T., Brown, P. O. and Botstein, D. (1998) Cluster analysis and display of genome-wide expression patterns. *Proceedings of the National Academy of Sciences*, **95**, 14863–14868.
- Fern, X. Z. and Brodley, C. E. (2003) Random projection for high dimensional data clustering: A cluster ensemble approach. In *Proceedings of the 20th International Conference on Machine Learning*, 186–193.
- Fraley, C. and Raftery, A. (1998) MCLUST: Software for model-based cluster and discriminant analysis. *Department of Statistics, University of Washington: Technical Report*, **342**, 1312.
- Gataric, M., Wang, T. and Samworth, R. J. (2020) Sparse principal component analysis via axis-aligned random projections. *J. Roy. Statist. Soc., Ser. B*, **82**, 329–359.
- Han, S. and Boutin, M. (2015) The hidden structure of image datasets. In *2015 IEEE International Conference on Image Processing (ICIP)*, 1095–1099, IEEE.
- Hastie, T., Tibshirani, R. and Friedman, J. H. (2009) *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, vol. 2. Springer.
- Ho, N., Khamaru, K., Dwivedi, R., Wainwright, M. J., Jordan, M. I. and Yu, B. (2020) Instability, computational efficiency and statistical accuracy. *arXiv preprint arXiv:2005.11411*.
- Ingster, Y. and Suslina, I. A. (2012) *Nonparametric Goodness-of-Fit Testing Under Gaussian Models*, vol. 169. Springer Science & Business Media.
- Jain, A. K. and Flynn, P. J. (1996) Image segmentation using clustering. In *Advances in image understanding: A Festschrift for Azriel Rosenfeld*, 65–83, IEEE Press, Piscataway, NJ.
- Jin, J. and Wang, W. (2016) Influential features PCA for high dimensional clustering. *Ann. Statist.*, **44**, 2323–2359.
- Johnson, W. B. and Lindenstrauss, J. (1984) Extensions of Lipschitz maps into a Hilbert space. *Contemp. Math.*, **26**, 189–206.
- Kaufman, L. and Rousseeuw, P. J. (2009) *Finding Groups in Data: an Introduction to Cluster Analysis*, vol. 344. John Wiley & Sons.
- Laurent, B. and Massart, P. (2000) Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.*, **28**, 1302–1338.
- Liang, P. (2005) *Semi-supervised learning for natural language*. Ph.D. thesis, Massachusetts Institute of Technology.
- Liu, H., Gao, C. and Samworth, R. J. (2021) Minimax rates in sparse, high-dimensional change point detection. *Ann. Statist.*, **49**, 1081–1112.
- Lloyd, S. (1982) Least squares quantization in PCM. *IEEE Transactions on Information Theory*, **28**, 129–137.

- Löffler, M., Wein, A. S. and Bandeira, A. S. (2022) Computationally efficient sparse clustering. *Information and Inference: A Journal of the IMA*, **11**, 1255–1286.
- Löffler, M., Zhang, A. Y. and Zhou, H. H. (2021) Optimality of spectral clustering in the Gaussian mixture model. *Ann. Statist.*, **49**, 2506–2530.
- Lopes, M., Jacob, L. and Wainwright, M. J. (2011) A more powerful two-sample test in high dimensions using random projection. In *Advances in Neural Information Processing Systems*, 1206–1214.
- Mai, Q., Zou, H. and Yuan, M. (2012) A direct approach to sparse discriminant analysis in ultra-high dimensions. *Biometrika*, **99**, 29–42.
- Marzetta, T. L., Tucci, G. H. and Simon, S. H. (2011) A random matrix-theoretic approach to handling singular covariance estimates. *IEEE Transactions on Information Theory*, **57**, 6256–6271.
- Muirhead, R. J. (2009) *Aspects of Multivariate Statistical Theory*. John Wiley & Sons.
- Ndaoud, M. (2022) Sharp optimal recovery in the two component Gaussian mixture model. *Ann. Statist.*, **50**, 2096–2126.
- Oymak, S. and Gulcu, T. C. (2020) Statistical and algorithmic insights for semi-supervised learning with self-training. *Preprint*, arxiv:2006.11006.
- Ramey, J. A. (2016) Datamicroarray: Collection of data sets for classification. *R Package*, <https://rdrr.io/github/ramhiser/datamicroarray/>.
- Reeve, H. W., Kabán, A. and Bootkrajang, J. (2022) Heterogeneous sets in dimensionality reduction and ensemble learning. *Machine Learning*, 1–22.
- Rodriguez, M. Z., Comin, C. H., Casanova, D., Bruno, O. M., Amancio, D. R., Costa, L. d. F. and Rodrigues, F. A. (2019) Clustering algorithms: a comparative approach. *PLoS ONE*, **14**, e0210236.
- Slawski, M. (2018) On principal components regression, random projections, and column subsampling. *Electron. J. Statist.*, **12**, 3673–3712.
- Stewart, G. W. and Sun, J. (1990) *Matrix Perturbation Theory*. Academic Press, Inc., San Diego, CA.
- Thanei, G.-A., Heinze, C. and Meinshausen, N. (2017) Random projections for large-scale regression. In *Big and Complex Data Analysis*, 51–68, Springer.
- Turian, J., Ratinov, L. and Bengio, Y. (2010) Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 384–394.
- Van Engelen, J. E. and Hoos, H. H. (2020) A survey on semi-supervised learning. *Machine Learning*, **109**, 373–440.
- Vershynin, R. (2012) How close is the sample covariance matrix to the actual covariance matrix? *Journal of Theoretical Probability*, **25**, 655–686.

- Verzelen, N. and Arias-Castro, E. (2017) Detection and feature selection in sparse mixture models. *Ann. Statist.*, **45**, 1920–1950.
- von Luxburg, U. (2007) A tutorial on spectral clustering. *Statistics and Computing*, **17**, 395–416.
- Wang, D., Lin, J., Cui, P., Jia, Q., Wang, Z., Fang, Y., Yu, Q., Zhou, J., Yang, S. and Qi, Y. (2019) A semi-supervised graph attentive network for financial fraud detection. In *2019 IEEE International Conference on Data Mining (ICDM)*, 598–607, IEEE.
- Wasserman, L., Azizyan, M. and Singh, A. (2014) Feature selection for high-dimensional clustering. *Preprint*, arxiv:1406.2240.
- Weyl, H. (1912) Das asymptotische Verteilungsgesetz der Eigenwerte linearer partieller Differentialgleichungen (mit einer Anwendung auf der Theorie der Hohlraumstrahlung). *Math. Ann.*, **71**, 441–479.
- Witten, D. M. and Tibshirani, R. (2010) A framework for feature selection in clustering. *J. Amer. Statist. Assoc.*, **105**, 713–726.
- Witten, D. M. and Tibshirani, R. (2011) Penalized classification using Fisher’s linear discriminant. *J. Roy. Statist. Soc., Ser. B*, **73**, 753–772.
- Wu, Y. and Zhou, H. H. (2022) Randomly initialised EM algorithm for two-component Gaussian mixture achieves near optimality in $O(\sqrt{n})$ iterations. *Mathematical Statistics and Learning*, **4**, 143–220.
- Xu, D. and Tian, Y. (2015) A comprehensive survey of clustering algorithms. *Annals of Data Science*, **2**, 165–193.
- Xu, J., Hsu, D. J. and Maleki, A. (2016) Global analysis of Expectation Maximization for mixtures of two Gaussians. *Adv. Neur. Inf. Proc. Syst.*, **29**.
- Xu, R. and Wunsch, D. (2005) Survey of clustering algorithms. *IEEE Transactions on neural networks*, **16**, 645–678.
- Yan, B., Yin, M. and Sarkar, P. (2017) Convergence of gradient EM on multi-component mixture of Gaussians. In *Adv. Neur. Inf. Proc. Syst.*, 6956–6966.
- Yang, F., Liu, S., Dobriban, E. and Woodruff, D. P. (2021) How to reduce dimension with PCA and random projections? *IEEE Transactions on Information Theory*, **67**, 8154–8189.
- Yellamraju, T. and Boutin, M. (2018) Clusterability and clustering of images and other “real” high-dimensional data. *IEEE Transactions on Image Processing*, **27**, 1927–1938.
- Zhang, A., Brown, L. D. and Cai, T. T. (2019) Semi-supervised inference: general theory and estimation of means. *Ann. Statist.*, **47**, 2538–2566.
- Zhu, X. and Goldberg, A. B. (2009) Introduction to semi-supervised learning. In *Synthesis Lectures on Artificial Intelligence and Machine Learning* (Brachman, R. J. and Dietterich, T., eds.), 1–130, Morgan & Claypool Publishers.
- Zhu, X. J. (2005) Semi-supervised learning literature survey. Tech. rep., University of Wisconsin-Madison Department of Computer Sciences.