LEARNING HIGH-DIMENSIONAL GAUSSIAN MIXTURE MODELS VIA A FOURIER APPROACH

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we address the challenge of learning high-dimensional Gaussian mixture models (GMMs), with a specific focus on estimating both the model order and the mixing distribution from i.i.d. samples. We propose a novel algorithm that achieves linear complexity relative to the sample size n, significantly improving computational efficiency. Unlike traditional methods, such as the method of moments or maximum likelihood estimation, our algorithm leverages Fourier measurements from the samples, facilitating simultaneous estimation of both the model order and the mixing distribution. The difficulty of the learning problem can be quantified by the separation distance Δ and minimal mixing weight w_{\min} . For stable estimation, a sample size of $\Omega\left(\frac{1}{w_{\min}^2\Delta^{4K-4}}\right)$ is required for the model order, while $\Omega\left(\frac{1}{w_{\min}^2\Delta^{4K-2}}\right)$ is necessary for the mixing distribution. This highlights the distinct sample complexities for the two tasks. For D-dimensional mixture models, we propose a PCA-based approach to reduce the dimension, reducing the algorithm's complexity to $O(nD^2)$, with potential further reductions through random projections. Numerical experiments demonstrate the efficiency and accuracy compared with the EM algorithm. In particular, we observe a clear phase transition in determining the model order, as our method outperforms traditional information criteria. Additionally, our framework is flexible and can be extended to learning mixtures of other distributions, such as Cauchy or exponential distributions.

031 032

033 034

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

1 INTRODUCTION

1.1 BACKGROUND

The Gaussian Mixture Model (GMM) is a widely used statistical model that has found numerous applications in various fields, including machine learning, pattern recognition, data clustering, and image processing. It is a powerful tool for modeling complex data and signals originating from sub-populations or distinct sources. The GMM represents a probability distribution as a weighted sum of Gaussian components, each characterized by its mean and covariance matrix. Formally, each observation of the GMM follows:

- 042
- 043
- 044 045

046

047

048

050 051

052

where w_i is the mixing weight such that $w_i > 0$ and $\sum_{i=1}^{K} w_i = 1$. The mean and the covariance matrix of the *i*-th component are denoted as μ_i and Σ_i , respectively. For each sample x, we can introduce a latent variable $z \in \{1, \dots, K\}$, with the marginal distribution of z specified by the mixing weights:

 $oldsymbol{x} \sim \sum_{i=1}^{K} w_i \mathcal{N}(oldsymbol{\mu}_i, oldsymbol{\Sigma}_i)$

$$\mathbb{P}(z=i)=w_i.$$

Thus, the GMM can also be expressed conditionally as

$$\boldsymbol{x}|(z=i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i).$$
 (2)

(1)

Given the i.i.d. samples drawn from the mixture distribution, the challenge is to learn the underlying model. Generally, there are three formulations for learning mixtures:

055 056

058

059 060

061 062

063 064

073

074

075

077 078 079

087

- *Clustering*: estimate the latent variable z_i for each sample x_i ;
- Parameter estimation: estimate the weights w_i 's, means μ_i 's and covariance matrix Σ_i 's up to a global permutation;
- *Density estimation*: estimate the probability density function of the GMM under specific loss functions.

Existing methodologies for clustering primarily rely on k-means, which seeks to minimize:

$$\arg\min_{z_j, \boldsymbol{\mu}_i} \sum_{j=1}^n \sum_{i=1}^K \mathbf{1} \{ z_j = i \} \| \boldsymbol{x}_j - \boldsymbol{\mu}_i \|^2,$$
(3)

065 where $1\{z_i = i\} = 1$ if $z_i = i$ otherwise 0. It is well-known that solving the k-means exactly 066 in the general case is NP-hard, even for two clusters (see Aloise et al. (2009)). Various computa-067 tionally tractable approximation approaches have been proposed, including the widely used Lloyd's algorithm (Lloyd (1982)), nonnegative matrix factorization (NMF) (see Paatero & Tapper (1994); 068 He et al. (2011); Zhuang et al. (2023)), and semidefinite programming (SDP) (see Peng & Wei 069 (2007)). Note that Lloyd's algorithm iterates a two-phase of re-assigning the samples to clusters and re-computing the cluster means until convergence. The perfect clustering of the mixture depends on 071 the separation distance defined as: 072

$$\Delta := \min_{1 \le i < j \le K} \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|.$$
(4)

It has been shown in Ndaoud (2022) that the critical threshold for a perfect clustering of a twocomponent Gaussian mixture with a unified covariance matrix $\sigma^2 I$ in p-dimension is: 076

$$\Delta^2 = \sigma^2 \left(1 + \sqrt{1 + \frac{2p}{n \log n}} \right) \log n \tag{5}$$

080 Similar results are obtained for the K-component mixture model in Chen & Yang (2021).

081 Parameter estimation and density estimation benefit from a larger sample size, contrasting with the 082 perfect clustering scenario (5). Existing methodologies for learning the mixture can be broadly 083 categorized into the maximum-likelihood method and the moment-based method. The maximum 084 likelihood method aims to maximize the likelihood of the given samples. The likelihood function is 085 defined as

$$L(\boldsymbol{x}_j$$
's $|w_i$'s, $\boldsymbol{\mu}_i$'s, $\boldsymbol{\Sigma}_i$'s $) = \prod_{j=1}^n \left(\sum_{i=1}^K w_i g(\boldsymbol{x}_j; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \right),$

where $q(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the probability density function of Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance Σ . Numerous iterative methods for optimization are proposed to seek the maximum or local maximum of the likelihood function. Among them, the most widely used one is the 091 EM(Expectation-Maximization) Algorithm(Dempster et al. (1977)). The EM algorithm iterates a 092 two-step operation to find a local maximum of the logarithm likelihood function, which may not necessarily be the ground-truth parameters. The Lloyd's algorithm can be regarded as a determinis-094 tic version of the EM algorithm. The moment-based methods date back to Pearson (1894). However, 095 Pearson's method has practical limitations due to its sensitivity to moment selection and the instabil-096 ity of finding roots of high-degree polynomials. Various modifications of the method of moments are proposed, such as the Generalized Method of Moments(Hansen (1982)) and the Denoised Method 098 of Moments(Wu & Yang (2020)). The Markov Chain Monte Carlo (MCMC) methods are also com-099 monly used to generate parameter samples from the posterior distribution, with prominent samplers including the Metropolis method (Metropolis et al. (1953)). Additionally, relating to this paper, 100 Fourier approach is proposed and utilized to learn the GMMs in Qiao et al. (2022); Liu & Zhang 101 (2024).102

103 It is worth noting that both the clustering via k-means and the parameter estimation by maximum 104 likelihood and moment-based methods require the model order K as an input. However, the model 105 order is often unknown a priori, necessitating a method for determining the appropriate order for model learning. To address the challenge of model order selection, various statistical criteria and 106 information-theoretic measures have been proposed. These include the Akaike Information Cri-107 terion (AIC, Akaike (1998)) and Bayesian Information Criterion (BIC, Schwarz (1978)). These methods aim to balance model complexity and goodness of fit, providing a quantitative measure to
evaluate the trade-off between model complexity and data fidelity. Bayesian approaches can also be
used to determine the model order. The variational inference method proposed by Corduneanu &
Bishop (2001) allows for model order determination by assigning appropriate prior distributions to
the parameters and maximizing the variational posterior distribution.

114 1.2 PROBLEM SETTING AND MAIN CONTRIBUTIONS

1

Given n independent samples drawn from a D-dimensional Gaussian mixture distribution with a unified covariance matrix:

$$\boldsymbol{x}_j \sim \sum_{i=1}^{K} w_i \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}), \quad j = 1, \cdots n.$$
 (6)

We assume that the covariance matrix $\Sigma \in \mathbb{R}^{D \times D}$ is known as prior information. This scenario is referred to as the Gaussian location mixture if $\Sigma = \sigma^2 I$. We define the separation distance Δ and the minimal weight w_{\min} of the model (6) as

$$\Delta = \min_{1 \le i < i \le K} \left\| \boldsymbol{\mu}_i - \boldsymbol{\mu}_j \right\|, \quad w_{\min} = \min_{1 \le i \le K} w_i.$$

In this paper, we focus on the parameter estimation of the mixture model (6). Specifically, we aim to determine the model order K and estimate the mixing distribution $\nu(x) = \sum_{i=1}^{K} w_i \delta_{\mu_i}(x)$ of the model from the independent samples.

130 **Contributions:** We propose an efficient algorithm to estimate the model order and parameters si-131 multaneously for high-dimensional GMMs, extending the previous work in Liu & Zhang (2024) for 132 one dimension. The time complexity of the proposed algorithm is linear in the sample size n, mak-133 ing it highly scalable. The main novelty of our approach is the leverage of the Fourier measurements 134 of the samples. This is naturally connected to the problem of super-resolution and of line spectral 135 estimation, which can be solved efficiently using subspace methods such as the MUltiple SIgnal 136 Classification (MUSIC) algorithm. To handle high-dimensional data, we apply Principal Compo-137 nent Analysis (PCA) to reduce the dimension to reduce the complexity to $O(D^2)$, significantly 138 improving computational efficiency in large-scale applications. We compare our algorithm with the 139 EM algorithm to highlight its advantages across different scenarios. We note that the Fourier approach in this paper differs from the one in Qiao et al. (2022), which primarily focuses on spherical 140 GMMs in low-dimensional settings and is based on estimating the Fourier transform of the mixture 141 at carefully chosen frequencies. 142

We establish a fundamental limit to estimating the model order and mixing distribution in the mixture model using the Fourier measurements. Specifically, we show that stable recovery of the model order requires a sample size of $n = \Omega\left(\frac{1}{w_{\min}^2 \Delta^{4K-4}}\right)$, while stable estimation of the means requires $n = \Omega\left(\frac{1}{w_{\min}^2 \Delta^{4K-2}}\right)$, respectively. This result quantifies the distinct sample complexities for these two tasks. We also provide multiple comparison tests with other model order estimation methods and illustrate a phase transition in the estimation accuracy.

150 151

152

113

122

123

124 125

126

1.3 PAPER ORGANIZATION AND NOTATIONS

The rest of the paper is organized as follows. In Section 2, we propose Algorithm 1 for model order
and mixing distribution estimation of GMMs and establish the sample size guarantee for stable estimation. In Section 3, we use PCA to reduce the time complexity of Algorithm 1 in high-dimensional
mixtures. We performed several numerical experiments to illustrate the accuracy, resolution, and efficiency of the algorithms in Section 2.

158 Throughout the paper, we write f(n) = O(g(n)) if there exists some constant $c_1 > 0$ such that 159 $f(n) < c_1g(n)$, and $f(n) = \Omega(g(n))$ if there exists some constant $c_2 > 0$ such that $f(n) > c_2g(n)$. 160 We denote $f(n) \simeq g(n)$ if $f(n) = \Omega(g(n))$ and f(n) = O(g(n)). For a k-dimensional subspace W161 of \mathbb{R}^n , the projection of vector $v \in \mathbb{R}^n$ on to W is defined as $\operatorname{Proj}_W(v) = \arg\min_{u \in W} ||u - v||_2$. 162 I_D denotes the identity matrix of rank D.

2 OUR PROPOSAL: MODEL ORDER AND MIXING DISTRIBUTION ESTIMATION VIA FOURIER APPROACH

2.1 Algorithm

162

163

164

166 167

171 172 173

174

178

179

181

182

183

185 186

187 188 189

190 191 192

193

In this section, we present our algorithm for model order and mixing distribution estimation of highdimensional GMMs. Our approach leverages the Fourier transform of the mixture distribution and highlights a natural connection with line spectral estimation (LSE) and super-resolution (SR). We assume that the means $\mu_i \in [-R, R)^D$ for some R > 0. The probability density function of the distribution (6) can be expressed in a convolutional form:

$$p(\boldsymbol{x}) = g(\boldsymbol{x}; \boldsymbol{\Sigma}) * \sum_{i=1}^{K} w_i \delta_{\boldsymbol{\mu}_i}(\boldsymbol{x})$$
(7)

where $g(\boldsymbol{x}; \boldsymbol{\Sigma})$ is the density function of the Gaussian distribution $\mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$. We now consider the Fourier transform of (7):

$$\phi(\boldsymbol{t}) = \mathcal{F}[p(\boldsymbol{x})] = e^{-\boldsymbol{t}^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{t}} \sum_{i=1}^{K} w_i e^{\iota \langle \boldsymbol{\mu}_i, \boldsymbol{t} \rangle}, \qquad (8)$$

where $\mathcal{F}[\cdot]$ denotes the Fourier transform. The function $\phi(t)$ is also known as the characteristic function (CF) of the mixture model in the context of the statistics. It can be estimated from the samples by the empirical characteristic function (ECF):

$$\psi_n(\boldsymbol{t}) = \frac{1}{n} \sum_{j=1}^n e^{\iota \langle \boldsymbol{x}_j, \boldsymbol{t} \rangle}.$$
(9)

According to the central limit theorem, the ECF follows asymptotic normality:

$$\sqrt{n} (\psi_n(\boldsymbol{t}) - \phi(\boldsymbol{t})) \xrightarrow{d} \mathcal{N}(0, 1 - |\phi(\boldsymbol{t})|^2), \quad n \to +\infty.$$

By modulating (9) with the term $e^{t^{T}\Sigma t}$, we obtain:

$$e^{\boldsymbol{t}^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{t}}\psi_{n}(\boldsymbol{t}) = \sum_{i=1}^{K} w_{i}e^{\iota\langle\boldsymbol{\mu}_{i},\boldsymbol{t}\rangle} + \epsilon_{n}(\boldsymbol{t}), \qquad (10)$$

194 where the right-hand side consists of a linear combination of exponential signals and a noise term 195 $\epsilon_n(t)$ that is due to the finite sample size n. The estimation of μ_i 's from the measurement (10) 196 is known as the Line Spectral Estimation (LSE), see Stoica et al. (2005). Due to the exponential 197 decay of the Fourier data $\phi(t)$, the available measurement in (10) is band-limited in the sense that there exist positive numbers f_1, f_2, \dots, f_D , called cutoff frequencies, such that only measurement at $t = (t_1, \dots, t_D)$ with $|t_i| \le f_i$ for $1 \le i \le D$ can be used for estimation. Estimating μ_i 's when 199 they are closely separated from the band-limited Fourier data is a super-resolution problem, see 200 Donoho (1992). The success of LSE depends crucially on the noise level and the cutoff frequencies. 201 In our problem, the noise level $\|\epsilon_n(t)\|_{\infty}$ can be estimated quantitatively in a probabilistic manner 202 by the following proposition: 203

Proposition 1. For any fixed $\epsilon > 0$, we have

$$\mathbb{P}\left(\left|e^{\mathbf{t}^{\mathrm{T}}\mathbf{\Sigma}\mathbf{t}}\psi_{n}(\mathbf{t})-\sum_{i=1}^{K}w_{i}e^{\iota\langle\boldsymbol{\mu}_{i},\mathbf{t}\rangle}\right|\geq\epsilon\right)\leq4\exp\left\{-\frac{n\epsilon^{2}}{4e^{2\mathbf{t}^{\mathrm{T}}\mathbf{\Sigma}\mathbf{t}}}\right\}\leq4\exp\left\{-\frac{n\epsilon^{2}}{4e^{2\|\mathbf{t}\|_{2}^{2}\sigma_{\max}(\mathbf{\Sigma})}}\right\}$$

where $\sigma_{\max}(\Sigma)$ denotes the maximal singular value of Σ . Then for any $\delta \in (0, 1)$, if the sample size $n \ge 4 \log\left(\frac{4}{\delta}\right) \frac{e^{2\|\mathbf{t}\|_2^2 \sigma_{\max}(\Sigma)}}{\epsilon^2}$, with probability $1 - \delta$, we have that

210 211

209

 $\epsilon_{e} = \epsilon_{e} = \epsilon_{e}$, with probability 1 = 0, we have ϵ_{e} , $\epsilon_{e} = \epsilon_{e}$.

Given measurement (10) at a uniform grid of domain $[-f_1, f_1] \times \cdots [-f_D, f_D]$, we employ a MUSIC-type algorithm to estimate both the model order K (i.e., the number of Gaussian components) and mixing distribution $\nu(\boldsymbol{x}) = \sum_{i=1}^{K} w_i \delta_{\mu_i}(\boldsymbol{x})$. Introduced by Schmidt (1986), the MUltiple SIgnal Classification (MUSIC) algorithm is a widely utilized technique in frequency estimation, 216 spectral analysis, and radar signal processing, renowned for its high-resolution parameter estimation 217 capabilities. Essentially, the MUSIC algorithm exploits the exponential form of the signals (similar 218 to Prony's method introduced in Prony (1795)), as defined in Equation (10), to construct a Hankel 219 matrix that admits a Vandermonde decomposition. The algorithm proceeds by performing Singular 220 Value Decomposition (SVD) on the Hankel matrix to identify the noise subspace. Subsequently, it formulates an imaging function (denoted as $\mathcal{J}(\mu)$ in the algorithm) by computing a noise-space 221 correlation function. In the noiseless scenario, the imaging function exhibits peaks precisely at the 222 set of Gaussian means $\{\mu_i\}_{1 \le i \le K}$. In the presence of noise, the algorithm determines the number of Gaussian means by identifying the number of local maxima in the imaging function and esti-224 mates the set of means based on the locations of these maxima. The details of the MUSIC algorithm 225 can be found in Appendix B. In Section 2.2 and 2.3, we discuss how to select the cutoff frequen-226 cies f_1, \dots, f_D and the number of sampling points to balance the computational tractability and 227 estimation accuracy. The mixing weights are estimated using the quadratic programming, as de-228 tailed in Appendix F. We summarize the model order selection and mixing distribution estimation 229 in Algorithm 1. 230

232

233 234

237

238

242 243 244

245

246 247

248

253

254 255

256

257 258

259 260 261

262

263 264

265

input : samples X_1, \dots, X_n , covariance matrix Σ , cutoff frequencies (f_1, \dots, f_D) , a prior upper bound for the number of Gaussian components L, sample size of the Fourier measurement in each direction N with N > L + K.

²³⁵ ²³⁶ 1 Compute $y_n(t) = e^{t^T \Sigma t} \psi_n(t)$ on the uniform grid of $[-f_1, f_1] \times \cdots \times [-f_D, f_D]$ with (N+1) sample points along each direction;

Algorithm 1: Model Order Selection and Mixing Distribution Estimation

2 Apply Algorithm 3 with input $y_n(t)$, N, L and plot the imaging function $\mathcal{J}(\mu)$ in $[-R, R)^D$;

3 Return the model order \hat{K} as the number of local maxima of $\mathcal{J}(\mu)$ and the means as the local maxima $\{\hat{\mu}_i\}_{1 \le i \le \hat{K}}$;

4 Return the weights $\{\hat{w}_i\}_{1 \le i \le \hat{K}}$ by solving the quadratic programming problem (35);

output: estimated mixing distribution $\hat{\nu}(\boldsymbol{x}) = \sum_{i=1}^{\hat{K}} \hat{w}_i \delta_{\hat{\mu}_i}(\boldsymbol{x}).$

We remark that this algorithm is also applicable when the model order K is known. In that case, the means are determined by selecting the largest K local maxima of the imaging function $\mathcal{J}(\mu)$.

2.2 PARAMETER SETUP OF ALGORITHM 1

In this section, we discuss how to set the parameters in Algorithm 1. Recall the Gaussian means μ_j 's are located within $[-R, R)^D$. By the the Nyquist–Shannon sampling theorem, the sampling step size h for each direction in the Fourier domain should satisfy $0 < h \le \frac{\pi}{R}$, resulting in the following condition on the sample size N in each direction:

$$N \ge \frac{2f_d R}{\pi}, \quad d = 1, \cdots, D.$$

We assume that we have a prior upper bound L of the number of Gaussian components K with L = O(K). To recover L components by the MUSIC algorithm, a sufficient condition on N (see Appendix B) is:

$$N \ge 2L + 1.$$

Therefore

$$N = \max\left(2L+1, \left\lceil \frac{2f_{\max}R}{\pi} \right\rceil\right),\tag{11}$$

where $f_{\max} = \max\{f_d : d = 1, \dots, D\}$ and $\lceil x \rceil$ is the smallest integer greater or equal to x. The choice of cutoff frequencies will be discussed in Section 2.3.

2.3 TIME AND SAMPLING COMPLEXITY OF ALGORITHM 1

In this section, we analyze the time and sampling complexity of Algorithm 1 with parameters set as (11). We also propose a method for determining the cutoff frequencies. The time complexity of computing the Fourier measurement $y_n(t)$ is given by:

$$O\left(n(N+1)^D\right)$$

For multidimensional MUSIC, the primary computational cost arises from the singular value decomposition and the evaluation of the imaging function $\mathcal{J}(\mu)$. Suppose the number of grid points for evaluating $\mathcal{J}(\mu)$ is *M* along each dimension. The time complexity of Algorithm 3 is:

$$O\left(\min\left\{(L+1)^{2D}(N-L+1)^{D},(L+1)^{D}(N-L+1)^{2D}\right\}+(2M)^{D}\right)$$

²⁷⁵ Using the inputs from (11), the overall time complexity of Algorithm 1 becomes

276 277

287 288 289

292

293 294 295

304

305

307 308

309

310311312313

321

322 323

273 274

$$O(n2^D K^D + K^{3D} + 2^D M^D)$$
(12)

which is linear in sample size n, but exponential in dimensionality D. This complexity can be reduced using the dimension reduction method introduced in Section 3.

Next, we examine the sampling complexity in relation to the separation distance Δ and the minimal mixing weight w_{\min} . The reliability of the estimation provided by Algorithm 1 depends on these two parameters, as well as the noise level $|\epsilon_n(t)|$ which is determined by the sample size n. This relationship is closely connected to the computational resolution limits established in Liu & Zhang (2021b) for one-dimensional and Liu & Zhang (2021a) for multi-dimensional LSE. Before presenting the main theorem, we introduce the concept of the computational resolution limit for multi-dimensional LSE. Consider the multi-dimensional Fourier measurement defined as

$$y(\boldsymbol{t}) = \sum_{i=1}^{K} w_i e^{\iota \langle \boldsymbol{\mu}_i, \boldsymbol{t} \rangle} + \epsilon(\boldsymbol{t}), \quad \boldsymbol{t} \in \mathbb{R}^D, \quad \|\boldsymbol{t}\|_2 \le f.$$
(13)

Assume that $\|\epsilon(t)\|_{\infty} < \sigma$.

Definition 1. Given the Fourier measurement y(t) in (13), we say that the $\hat{\nu}(\boldsymbol{x}) = \sum_{i=1}^{\hat{K}} \hat{w}_i \delta_{\hat{\mu}_i}(\boldsymbol{x})$ is a σ -admissible discrete measure of y(t) if

$$\left\|\mathcal{F}\hat{\nu}(t) - y(t)\right\|_{\infty} < \sigma, \quad \forall \left\|t\right\|_{2} \le f_{*}$$

The set of σ -admissible measures characterizes all the possible solutions of the inverse problem from Fourier measurements y(t). If there exists an admissible measure $\hat{\nu}$ with less than K components, one may miss out one or more sources and therefore cannot estimate the model order correctly. This leads to the definition of the computational resolution limit for number detection.

Definition 2. The computational resolution limit for number detection in D-dimensional space is defined as the smallest nonnegative number $\mathcal{R}_{D,K}$ such that for all K-component measure $\sum_{i=1}^{K} w_i \delta_{\mu_i}, \mu_i \in B_{\frac{(K-1)\pi}{2f}}^{D}(\mathbf{0})$ and the associated Fourier measurement $y(\mathbf{t})$ in (13), if

$$\Delta = \min_{1 \le i < j \le K} \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\| \ge \mathcal{R}_{D,K}$$

then there exists no σ -admissible measure consisting less than K components with Fourier measurements y(t).

A quantitative characterization of $\mathcal{R}_{D,K}$ is provided in Appendix D. It can be shown that, up to two constants depending only on D, K, the limit takes the form

$$\mathcal{R}_{D,K} \asymp \frac{\pi}{f} \left(\frac{\sigma}{w_{\min}}\right)^{\frac{1}{2K-2}}.$$
(14)

This computational limit indicates that to accurately estimate the model order from the Fourier measurements (10) of the samples. The noise level $\|\epsilon_n(t)\|$ must be small enough such that $\mathcal{R}_{D,K} \leq \Delta$. The following theorem establishes a non-asymptotic lower bound for the sample size required to accurately recover the model order.

Theorem 1. Consider the D-dimensional mixture model $\sum_{i=1}^{K} w_i \mathcal{N}(\mu_i, \Sigma)$ with $\mu_i \in B_{\frac{(K-1)\pi}{2f}}^D(\mathbf{0})$.

For any $\delta \in (0, 1)$ *, if the sample size* n *satisfies that*

$$n \ge C_{K,D} \log\left(\frac{4}{\delta}\right) \frac{e^{2f^2 \sigma_{\max}(\mathbf{\Sigma})}}{w_{\min}^2 (f\Delta)^{4K-4}}.$$
(15)

Then with probability $1 - \delta$, $\Delta \geq \mathcal{R}_{D,K}$ holds. Here $C_{K,D}$ is a constant only relying on K and D.

Remark 1. For exact estimation of the model order K using the Fourier measurements (10), the sample size should satisfy that

$$n = \Omega\left(\frac{1}{w_{\min}^2 \Delta^{4K-4}}\right). \tag{16}$$

This reveals the relation of the sample size n with the mixture model itself explicitly.

Remark 2. The computational resolution limit for support recovery has also been established in Liu & Zhang (2021a). Following this theory, the sample complexity for estimating the means of a K-component GMMs with an error threshold less than $\Delta/2$, where Δ is the separation distance of the means, should satisfy

$$n = \Omega\left(\frac{1}{w_{\min}^2 \Delta^{4K-2}}\right)$$

The computational resolution limit also sheds light on setting the cutoff frequencies in Algorithm 1. From Proposition 1, the noise level in (10) is amplified by a factor of $e^{t^T \Sigma t}$. For a one-dimensional mixture with variance σ^2 , the noise level is amplified by $e^{f^2 \sigma^2}$. To minimize the computational resolution limit, a straightforward calculation leads to the optimal cutoff frequency set as $f^{\text{optimal}} = \sqrt{1-2}$

$$\frac{2K-2}{\sigma^2}$$
. Therefore, we can set the cutoff frequencies as

V

 $f_d = \sqrt{\frac{2L-2}{\boldsymbol{e}_d^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{e}_d}}, \quad d = 1, \cdots, D,$ (17)

if K is unknown. Along with (11), these parameters are tested in detail in the numerical experiments shown in Section 4.

3 PCA-BASED DIMENSION REDUCTION

The time complexity of Algorithm 1 is exponential in the data dimension D (see (12)). For a K-component mixture model, the means S lies on a subspace at most dimension K. If we can identify this subspace and project the samples onto it, the computational complexity of the model order estimation can be significantly reduced. In this section, we introduce a PCA-based method for dimension reduction. The idea is to first project the data onto a low-dimensional linear manifold using Principle Component Analysis (PCA) before running Algorithm 1. We demonstrate that this projection-based technique can also be used to estimate the mixing distribution. The PCA is based on the Singular Value Decomposition(SVD) of the data matrix:

 $oldsymbol{X} = \left[oldsymbol{x}_1 \quad \cdots \quad oldsymbol{x}_n
ight]^{\mathrm{T}} \in \mathbb{R}^{n imes D}.$

Assume that n > D and denote its singular value decomposition as

$$oldsymbol{X} = \sum_{d=1}^D \lambda_d oldsymbol{u}_d oldsymbol{v}_d^{ ext{T}}$$

where $\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_D \ge 0$. Denote $V_k = [v_1 \cdots v_k] \in \mathbb{R}^{D \times k}$. The PCA projects the samples onto the column space of V_k . We summarize the PCA-based model order and mixing distribution estimation algorithm below.

8	Algorithm 2: PCA-based Model Order Selection and Mixture Distribution Estimation
9	input : samples $\boldsymbol{x}_1, \cdots, \boldsymbol{x}_n, \boldsymbol{\Sigma}, (f_1, \cdots, f_k), k, N, L$
0 1 2	1 Compute the SVD of data matrix $X = \sum_{d=1}^{D} \lambda_d \boldsymbol{u}_d \boldsymbol{v}_d^{\mathrm{T}}$; 2 Project the samples to the subspace spanned by $\boldsymbol{v}_1, \cdots, \boldsymbol{v}_k$;
2 3	3 Run Algorithm 1 with inputs as projected samples, $V^T \Sigma V$, (f_1, \dots, f_k) , N, L ; 4 Transfer the projected means into their original space;
4	output: the model order \hat{K} and the mixing distribution.

If the model order K is known, set k = K. If k < K, the projection may miss components lying orthogonal to the space spanned by the principle components. This issue and possible solutions are

discussed in Section C.2. In the next subsection, we provide a theoretical analysis of the dimension
 reduction using PCA. This method reduces the time complexity from (12) to

$$O\left(nD^2 + n2^kK^k + K^{3k} + 2^kM^k\right).$$
(18)

This reduces the exponential dependency on D to k. If K = O(1) relative to the sample size n and dimensionality D, then the time complexity of Algorithm 2 becomes $O(nD^2)$, which is quadratic in the dimensionality. This complexity can be further reduced by random projection techniques. For instance, one can first apply the Johnson-Lindestrauss embedding to project the data onto a subspace of dimension $\Omega\left(\frac{\log K}{\epsilon^2}\right)$. The estimation accuracy remains promising if the shrunk separation distance $(1 - \epsilon)\Delta$ remains above the resolution limit.

3.1 ANALYSIS ON THE GAUSSIAN LOCATION MIXTURE

In this section, we consider the Gaussian mixture with covariance matrix as $\sigma^2 I$, also known as the Gaussian location mixture. We demonstrate that when n > D, the expected subspace spanned by the first K right singular vectors $\{v_1, \dots v_K\}$ in PCA either includes or coincides with the subspace spanned by the means $\{\mu_1, \dots \mu_K\}$.

Firstly, recall that

$$\operatorname{span}\{\boldsymbol{v}_1,\cdots,\boldsymbol{v}_K\} = \operatorname*{arg\,max}_{\{\boldsymbol{V}:\dim\,\boldsymbol{V}=K\}} \|\operatorname{Proj}_{\boldsymbol{V}}\boldsymbol{X}\|_2.$$
(19)

We have the following theorem, where part of the proof is adapted from Vempala & Wang (2002):

Theorem 2. Given any k-dimensional $(k \leq D)$ subspace spanned by orthonormal vectors $\{w_1, \dots, w_k\}$. Denote $W_k = \operatorname{span}\{w_1, \dots, w_k\}$ and $U_K = \operatorname{span}\{\mu_1, \dots, \mu_K\}$, then we have

 $\mathbb{E} \left\| \operatorname{Proj}_{U_{K}} X \right\|_{2} \geq \mathbb{E} \left\| \operatorname{Proj}_{W_{k}} X \right\|_{2}$

Furthermore, if k < K, we have $\arg \max_{W_k} \mathbb{E} \| \operatorname{Proj}_{W_k} X \|_2 \subset U_K$ and if $k \geq K$, we have $U_K \subset \arg \max_{W_k} \mathbb{E} \| \operatorname{Proj}_{W_k} X \|_2$.

This theorem implies that, with a sufficiently large sample size, the subspace obtained via SVD closely approximates the subspace spanned by the centers. Therefore, estimating the mixing distribution in the projected subspace is reasonable.

- 412 4 NUMERICAL RESULTS
- 414 4.1 COMPARISON WITH EM ALGORITHM

In this experiment, we compare the performance of Algorithm 2 with the EM algorithm for estimating the mixing distribution. We also include tests using PCA as a preprocessing step before applying the EM algorithm. All tests are conducted on a mixture model with dimension D = 100and components K = 2, 3 with $\Sigma = I$.

The tests are designed as follows. For K = 2, samples are generated from the model 420 $\frac{1}{2}\mathcal{N}(-\mu, I_{100}) + \frac{1}{2}\mathcal{N}(\mu, I_{100})$, and the mixture distribution is $\frac{1}{3}\mathcal{N}(-\mu, I_{100}) + \frac{1}{3}\mathcal{N}(0, I_{100}) + \frac{1}{3}\mathcal{N}(0, I_{100})$ 421 $\frac{1}{3}\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{I}_{100})$ for $\tilde{K} = 3$. In each test, the model order K and the $\|\boldsymbol{\mu}\|_2$ are fixed. The sample size n ranges from 10,000 to 200,000 with increments of 10,000. With fixed $\|\boldsymbol{\mu}\|_2$, For 422 423 each sample size, 96 independent trials are conducted, and the estimation error is averaged across 424 trials. In each independent trial, the mean μ is generated by first selecting a vector uniformly 425 from the unit sphere \mathbb{S}^{D-1} , then scaling it by $\|\mu\|_2$. The inputs for Algorithm 2 are set as 426 $f = \sqrt{2K-2}, k = K, L = K, N = 2K$ in accordance with (11) and (17). For the EM algo-427 rithm, the initial means are randomly set as K samples, and the algorithm terminates after 5,000428 iterations or when the log-likelihood increases less than 1×10^{-6} . During the iterations of the EM 429 algorithm, the covariance matrix is fixed as I_{100} . The estimation error for the mixing distribution is 430 defined using the Wasserstein distance: 431

$$W_1(\nu, \hat{\nu}) = \inf \mathbb{E} \|X - Y\|_2$$

381

389

390

396 397 398

399

403 404

407

408

409

410 411



Figure 1: Comparison with the EM algorithm. The uppers are the W_1 errors of each method and the lowers are the average running time(seconds) of each trial. The samples of each trial comes from: (a) $\frac{1}{2}\mathcal{N}(-\mu, \mathbf{I}_{100}) + \frac{1}{2}\mathcal{N}(\mu, \mathbf{I}_{100}), \|\boldsymbol{\mu}\|_2 = 1$; (b) $\frac{1}{2}\mathcal{N}(-\mu, \mathbf{I}_{100}) + \frac{1}{2}\mathcal{N}(\mu, \mathbf{I}_{100}), \|\boldsymbol{\mu}\|_2 = 2$; (c) $\frac{1}{3}\mathcal{N}(-\mu, \mathbf{I}_{100}) + \frac{1}{2}\mathcal{N}(\mathbf{0}, \mathbf{I}_{100}) + \frac{1}{3}\mathcal{N}(\mu, \mathbf{I}_{100}), \|\boldsymbol{\mu}\|_2 = 1$; (d) $\frac{1}{3}\mathcal{N}(-\mu, \mathbf{I}_{100}) + \frac{1}{2}\mathcal{N}(\mathbf{0}, \mathbf{I}_{100}) + \frac{1}{3}\mathcal{N}(\mu, \mathbf{I}_{100}), \|\boldsymbol{\mu}\|_2 = 2$.

where the infimum is taken for all joint distributions of random vectors (X, Y) with marginals $\nu, \hat{\nu}$ and this Wasserstein distance can be numerically computed through optimal transport¹. The results are presented in Figure 1.

The results demonstrate that Algorithm 2 achieves comparable accuracy to the EM algorithm while 460 requiring significantly less time, especially for the large sample sizes. This efficiency arises because 461 the running time of Algorithm 2 scales linearly with the sample size n. In contrast, the EM algorithm 462 is highly sensitive to the initialization of the means and the landscape of the likelihood function, 463 which may result in slow convergence if the initialization or landscape is unfavorable. Additionally, 464 the time complexity of the EM algorithm in each iteration is $O(nKD^2)$, which is approximately the 465 total complexity order (18). The dependence of the convergence speed of the EM algorithm on the 466 likelihood function's landscape is evident when comparing panels (a) with (b) and (c) with (d) in 467 Figure 1. A larger separation distance results in a better landscape, leading to faster convergence for 468 the EM algorithm.

469 470 471

449

450

451

452 453

454 455 456

457

458

459

4.2 Resolution Limit and Phase Transition of Model Order Estimation

In this experiment, we explore the resolution limit of Algorithm 1 and compare it with other commonly used model order estimation methods. Specifically, we test the resolution limit for equally weighted two-component and four-component mixture model in \mathbb{R}^2 . The covariance matrices are *I* for all Gaussian components. The geometry of the component means is illustrated in Figure 2.

The tests are designed as follows. We uniformly take 2,800 $(\log_{10}(n), \Delta)$ points in the domain 477 $[2.5, 6.0] \times [0.2, 6.0]$. For each $(\log_{10}(n), \Delta)$ pair, we construct the equally weighted mixture model 478 with the means illustrated in Figure 2. We draw n independent samples from the model and apply 479 Algorithm 1, AIC, and BIC for model order estimation. For the K-component mixture, the inputs of 480 Algorithm 1 are $f = \sqrt{K+1}$, L = K+1, N = 2K+2, which allows for the model order ranging 481 from 1 to K + 1. For AIC and BIC, the model is estimated by the EM algorithm with model order 482 ranging from 1 to K + 1. The EM algorithm terminates after 5,000 iterations or the log likelihood 483 increases less than 1×10^{-5} . The results are shown in Figure 2. 484

¹In our experiments, we use wasserstein_distance_nd in the Python package scipy.

The results reveal phase transitions for all three methods. The proposed method demonstrates a
 more favorable phase transition region compared to the other two criteria. However, the transition
 may not be as pronounced as that of the information criteria. Further refinement of these criteria
 could enhance the performance of the proposed method.



Figure 2: Geometry of the means and the phase transition of different model order estimation methods. The black dots stand for the mean locations and Δ stands for the separation distance. The blue triangle means the model order is underestimated and the green triangle means the model order is overestimated.

515 516 517

518 519

513

514

5 CONCLUSIONS AND DISCUSSIONS

Learning Gaussian mixture models is a challenging task, particularly in high dimensions or when the 520 number of components is large or unknown. The performance of the learning algorithms depends 521 on the separation distance and minimal weight of the components. In this paper, we proposed an 522 efficient algorithm for estimating the model order and mixing distribution of the high-dimensional 523 GMMs. Our algorithm leverages the Fourier measurement of the samples, drawing a natural connec-524 tion to line spectral estimation and super-resolution techniques. We have established the sampling 525 complexities for estimating the model order and mixing distributions in relation to the separation 526 distance and minimal weight. Additionally, we demonstrated that the computational complexity for 527 learning high-dimensional mixtures can be further reduced using dimension reduction techniques 528 such as PCA. Empirical results confirmed that our algorithm achieves efficiency and accuracy comparable to, or better than, the EM algorithm. 529

530 We also acknowledge some aspects of our approach that present opportunities for future improve-531 ment. While our algorithm assumes that the unified covariance matrix Σ is known a priori, there are 532 scenarios where this may not be the case. To enhance the versatility of our method, estimating the 533 covariance matrix using Fourier measurements, as explored in the 1-D algorithm in Liu & Zhang 534 (2024), could be a promising direction. Additionally, while the time complexity remains quadratic with respect to dimensionality, this opens avenues for further research. Employing random pro-535 jections that preserve pairwise distances between components, such as the Johnson-Lindenstrauss 536 embedding (see Sanjeev & Kannan (2001)), could be an effective way to address this challenge. We 537 will exploring these possibilities in future work. 538

540 REFERENCES 541

554

558

559

561

562

565

566

567

568 569

570

571

572

573

576

581

- Hirotogu Akaike. Information theory and an extension of the maximum likelihood principle. In 542 Selected papers of hirotugu akaike, pp. 199–213. Springer, 1998. 543
- 544 Daniel Aloise, Amit Deshpande, Pierre Hansen, and Preyas Popat. Np-hardness of euclidean sumof-squares clustering. Machine learning, 75:245-248, 2009. 546
- Stephen Boyd and Lieven Vandenberghe. Convex optimization. Cambridge university press, 2004. 547
- 548 Xiaohui Chen and Yun Yang. Cutoff for exact recovery of gaussian mixture models. IEEE Transac-549 tions on Information Theory, 67(6):4223-4238, 2021. 550
- 551 Adrian Corduneanu and Christopher M Bishop. Variational bayesian model selection for mixture 552 distributions. In Artificial intelligence and Statistics, volume 2001, pp. 27–34. Morgan Kaufmann Waltham, MA, 2001. 553
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data 555 via the em algorithm. Journal of the royal statistical society: series B (methodological), 39(1): 556 1-22, 1977.
 - David L Donoho. Superresolution via sparsity constraints. SIAM journal on mathematical analysis, 23(5):1309–1331, 1992.
 - Zetao Fei and Hai Zhang. Scan-music: An efficient super-resolution algorithm for single snapshot wide-band line spectral estimation. arXiv preprint arXiv:2310.17988, 2023.
- 563 Lars Peter Hansen. Large sample properties of generalized method of moments estimators. Econo*metrica: Journal of the econometric society*, pp. 1029–1054, 1982. 564
 - Zhaoshui He, Shengli Xie, Rafal Zdunek, Guoxu Zhou, and Andrzej Cichocki. Symmetric nonnegative matrix factorization: Algorithms and applications to probabilistic clustering. IEEE Transactions on Neural Networks, 22(12):2117-2131, 2011.
 - Wenjing Liao. Music for multidimensional spectral estimation: Stability and super-resolution. IEEE transactions on signal processing, 63(23):6395–6406, 2015.
 - Ping Liu and Hai Zhang. A mathematical theory of computational resolution limit in multidimensional spaces. Inverse Problems, 37(10):104001, 2021a.
- 574 Ping Liu and Hai Zhang. A theory of computational resolution limit for line spectral estimation. 575 *IEEE Transactions on Information Theory*, 67(7):4812–4827, 2021b.
- Xinyu Liu and Hai Zhang. A fourier approach to the parameter estimation problem for one-577 dimensional gaussian mixture models. arXiv preprint arXiv:2404.12613, 2024. 578
- 579 Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2): 580 129–137, 1982.
- Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and Edward 582 Teller. Equation of state calculations by fast computing machines. The journal of chemical physics, 21(6):1087-1092, 1953. 584
- 585 Mohamed Ndaoud. Sharp optimal recovery in the two component gaussian mixture model. The 586 Annals of Statistics, 50(4):2096–2126, 2022.
- Pentti Paatero and Unto Tapper. Positive matrix factorization: A non-negative factor model with 588 optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126, 1994. 589
- Karl Pearson. Contributions to the mathematical theory of evolution. Philosophical Transactions of the Royal Society of London. A, 185:71–110, 1894. 592
- Jiming Peng and Yu Wei. Approximating k-means-type clustering via semidefinite programming. SIAM journal on optimization, 18(1):186-205, 2007.

594 595	R Prony. Essai experimental-, J. de l'Ecole Polytechnique, 2:929, 1795.
595 596	Mingda Qiao, Guru Guruganesh, Ankit Rawat, Kumar Ayinaya Dubey, and Manzil Zaheer. A fourier
597	approach to mixture learning. Advances in Neural Information Processing Systems, 35:20850–
598	20861, 2022.
599	
600 601	thirty-third annual ACM symposium on Theory of computing, pp. 247–257, 2001.
602 603	Tapan K Sarkar and Odilon Pereira. Using the matrix pencil method to estimate the parameters of a sum of complex exponentials. <i>IEEE Antennas and Propagation Magazine</i> , 37(1):48–55, 1995.
604 605 606	Ralph Schmidt. Multiple emitter location and signal parameter estimation. <i>IEEE transactions on antennas and propagation</i> , 34(3):276–280, 1986.
607	Gideon Schwarz. Estimating the dimension of a model. The annals of statistics, pp. 461–464, 1978.
608 609 610	Petre Stoica, Randolph L Moses, et al. <i>Spectral analysis of signals</i> , volume 452. Pearson Prentice Hall Upper Saddle River, NJ, 2005.
611 612	Gongguo Tang, Badri Narayan Bhaskar, and Benjamin Recht. Near minimax line spectral estimation. <i>IEEE Transactions on Information Theory</i> , 61(1):499–512, 2014.
613 614	Santosh Vempala and Grant Wang. A spectral algorithm for learning mixtures of distributions. In <i>The 43rd Annual IEEE Symposium on Foundations of Computer Science, 2002. Proceedings.</i> , pp.
615	113–122. IEEE, 2002.
615	Vibong Wu and Pengkun Yang. Optimal estimation of gaussian mixtures via denoised method of
618	moments. Annals of Statistics, 48(4), 2020.
619 620 621	Yubo Zhuang, Xiaohui Chen, Yun Yang, and Richard Y Zhang. Statistically optimal k-means clustering via nonnegative low-rank semidefinite programming. <i>arXiv preprint arXiv:2305.18436</i> , 2023.
622	
623	
624	
625	
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	
638	
639	
640	
641	
642	
643	
644	
645	
040 647	
047	

648 A NOTATIONS

We shall use the following notations in the appendix. We use $\Re(\cdot)$ and $\Im(\cdot)$ to denote the real and imaginary part of a complex number, vector or matrix. For a matrix $A \in \mathbb{R}^{m \times n}$ (or $\mathbb{C}^{m \times n}$), we use A_j to denote the *j*-th column of A and the induced 2-norm $\|A\|_2 = \sqrt{\sum_{j=1}^n \|A_j\|_2^2}$.

B REVIEWS ON MUSIC ALGORITHM

In this section, we review the multidimensional MUltiple SIgnal Classification(MUSIC) algorithm. The MUSIC algorithm (see Schmidt (1986)) was initially proposed for the direction of arrival(DOA) detection and line spectral estimation(LSE). The multidimensional MUSIC is applied in *D*-dimensional single-snapshot spectral estimation. Consider a signal y(t) which is a linear combination of K time-harmonic components and additive noise $\epsilon(t)$:

$$y(\boldsymbol{t}) = \sum_{i=1}^{K} w_i e^{\iota \langle \boldsymbol{\mu}_i, \boldsymbol{t} \rangle} + \epsilon(\boldsymbol{t}).$$
(20)

The goal is to recover the frequency set $S = \{\mu_i : 1 \le i \le K\}$ and the corresponding amplitude w_i , from uniform samples of y(t) in the domain $[-f, f]^D$. Suppose we have a total $(N + 1)^D$ uniformly spaced sampling points with a grid size $h = \frac{2f}{N}$. Consequently, the frequencies can only be determined on the torus $\left[-\frac{N\pi}{2f}, \frac{N\pi}{2f}\right]^D$.

670 We first review the multidimensional MUSIC algorithm when D = 2. The extension to higher 671 dimensions can be found in Liao (2015). For simplicity, we define the sampling coordinates along 672 each direction as $t_q = -f + q \frac{2f}{N}$ for $q = 0, \dots, N$, and $\mu_i = (\mu_1^i, \mu_2^i)$ for $i = 1, \dots, K$. We also 673 introduce the following notation:

$$\boldsymbol{\phi}_{l}(\boldsymbol{\mu}) = \begin{bmatrix} 1 & e^{\iota \boldsymbol{\mu} h} & \cdots & e^{\iota \boldsymbol{\mu} l h} \end{bmatrix}^{\mathrm{T}} \in \mathbb{C}^{l+1}.$$

677 Denote the noiseless uniform samples on the grid as:

$$u_{n_1,n_2} = \sum_{i=1}^{K} w_i e^{\iota \langle \boldsymbol{\mu}_i, \boldsymbol{t}_{n_1,n_2} \rangle}, \quad 0 \le n_1, n_2 \le N.$$

681 where $t_{n_1,n_2} = (t_{n_1}, t_{n_2})$ is the sample point.

Given a fixed integer L < N, we form the Hankel matrix

y

$$\boldsymbol{A}_{n_{1}} = \begin{bmatrix} y_{n_{1},0} & y_{n_{1},1} & \cdots & y_{n_{1},N-L} \\ y_{n_{1},1} & y_{n_{2},2} & \cdots & y_{n_{1},N-L+1} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n_{1},L} & y_{n_{1},L+1} & \cdots & y_{n_{1},N} \end{bmatrix} \in \mathbb{C}^{(L+1)\times(N-L+1)}, \quad 0 \le n_{1} \le N.$$
(21)

It is well known that A_{n_1} has the Vandermonde decomposition:

$$\boldsymbol{A}_{n_1} = \boldsymbol{\Phi}_{L,2} \boldsymbol{\Pi} \boldsymbol{\Lambda}_{n_1,1} \boldsymbol{\Phi}_{N-L,2}^{\mathrm{T}}, \qquad (22)$$

where

$$\begin{split} \Phi_{L,2} &= \begin{bmatrix} \phi_L(\mu_2^1) & \phi_L(\mu_2^2) & \cdots & \phi_L(\mu_2^K) \end{bmatrix} \in \mathbb{C}^{(L+1) \times K}, \\ \Pi &= \text{diag}(w_1, w_2 \cdots, w_K), \\ \Lambda_{n_1,1} &= \text{diag}\left(e^{\iota \mu_1^1 t_{n_1}}, e^{\iota \mu_1^2 t_{n_1}}, \cdots, e^{\iota \mu_1^K t_{n_1}}\right). \end{split}$$

⁶⁹⁷ Next, we construct the 2-fold Hankel block matrix:

For higher dimensions D > 2, the D-fold Hankel block matrix can be formed recursively as:

	$\begin{bmatrix} A_0 \end{bmatrix}$	A_1	•••	$A_{N_1-L_1}$	
H =	$\begin{vmatrix} A_1 \\ \vdots \end{vmatrix}$	$egin{array}{c} egin{array}{c} egin{array}$	•.	$\mathbf{A}_{N_1-L_1+1}$:	,
	\dot{A}_{L_1}	\dot{A}_{L_1+1}	•	$\dot{oldsymbol{A}}_{N_1}$	

where A_l is the (D - 1)-fold Hankel block matrix formed from the samples $\{y_{l,n_2,\cdots,n_D}: 0 \le n_2,\cdots,n_D \le N\}.$

For the 2-fold Hankel block matrix H, it can be verified that

$$\boldsymbol{H} = \begin{bmatrix} \boldsymbol{\Phi}_{L,2}\boldsymbol{\Lambda}_{0,1} \\ \boldsymbol{\Phi}_{L,2}\boldsymbol{\Lambda}_{1,1} \\ \vdots \\ \boldsymbol{\Phi}_{L,2}\boldsymbol{\Lambda}_{L,1} \end{bmatrix} \boldsymbol{\Pi} \begin{bmatrix} \boldsymbol{\Lambda}_{0,1}\boldsymbol{\Phi}_{N-L,2}^{\mathrm{T}} & \boldsymbol{\Lambda}_{1,1}\boldsymbol{\Phi}_{N-L,2}^{\mathrm{T}} & \cdots & \boldsymbol{\Lambda}_{N-L,1}\boldsymbol{\Phi}_{N-L,2}^{\mathrm{T}} \end{bmatrix}.$$
(24)

Defining:

$$\psi_L(\boldsymbol{\mu}) = \operatorname{vect}\left(\left\{e^{\iota\langle \boldsymbol{\mu}, \boldsymbol{t}_{n_1, n_2}\rangle} : 0 \le n_1, n_2 \le L\right\}\right) \in \mathbb{C}^{(L+1)^2},$$

the Vandermonde decomposition (24) can be written as:

$$\boldsymbol{H} = \underbrace{[\boldsymbol{\psi}_{L}(\boldsymbol{\mu}_{1}) \quad \cdots \quad \boldsymbol{\psi}_{L}(\boldsymbol{\mu}_{K})]}_{\boldsymbol{\Psi}_{L}} \boldsymbol{\Pi} \begin{bmatrix} \boldsymbol{\psi}_{N-L}(\boldsymbol{\mu}_{1}) \quad \cdots \quad \boldsymbol{\psi}_{N-L}(\boldsymbol{\mu}_{K}) \end{bmatrix}^{\mathrm{T}}.$$
 (25)

In the noiseless case, we have the following result for recovering the frequencies (see Liao (2015)): **Theorem 3.** Suppose $\mu_i \neq \mu_j$ for all $1 \leq i \neq j \leq K$ and

$$L+1 \ge K, \quad N-L+1 \ge K. \tag{26}$$

Then we have rank $(\Phi_{L,2}) = rank (\Phi_{N-L,2}) = rank (H) = K$. Furthermore, for any $\mu \in$ $\left[-\frac{Nw}{2f},\frac{Nw}{2f}\right]^2$, if (26) holds, we have

$$\boldsymbol{\mu} \in \mathcal{S} \Longleftrightarrow \boldsymbol{\psi}_L(\boldsymbol{\mu}) \in Im(\boldsymbol{\Psi}_L), \tag{27}$$

where $Im(\Psi_L)$ is the column space of Ψ_L .

This theorem provides a criterion (27) for detecting the frequencies in the noiseless case. When the measurement is contaminated with noise $\epsilon(t)$, we can apply the MUSIC algorithm by performing Singular Value Decomposition(SVD) on H^{ϵ} :

$$\boldsymbol{H}^{\epsilon} = \begin{bmatrix} \boldsymbol{U}_{1}^{\epsilon} & \boldsymbol{U}_{2}^{\epsilon} \end{bmatrix} \operatorname{diag} \begin{pmatrix} \sigma_{1}^{\epsilon}, \cdots, \sigma_{K}^{\epsilon}, \cdots \end{pmatrix} \begin{bmatrix} \boldsymbol{V}_{1}^{\epsilon} & \boldsymbol{V}_{2}^{\epsilon} \end{bmatrix}^{*},$$
(28)

where $U_1^{\epsilon} \in \mathbb{C}^{(L+1)^2 \times K}, U_2 \in \mathbb{C}^{(L+1)^2 \times \min\{(L+1)^2, (N-L+1)^2\} - K}$ and $Im(U_1^{\epsilon}), Im(U_2^{\epsilon})$ are called signal space and noise space, respectively. The algorithm is realized by projecting $\psi_L(\mu)$ onto the noise space and drawing the MUSIC imaging function defined as:

$$\mathcal{J}(\boldsymbol{\mu}) = \frac{\|\boldsymbol{\psi}_L(\boldsymbol{\mu})\|_2}{\|\boldsymbol{U}_2^{\epsilon*}\boldsymbol{\psi}_L(\boldsymbol{\mu})\|_2}.$$
(29)

In the noiseless case, we have the relation that

$$\boldsymbol{\mu} \in \mathcal{S} \iff \mathcal{J}(\boldsymbol{\mu}) = \infty.$$

In the noisy case, the frequency set S is determined by locating the local maxima of the imaging function $\mathcal{J}(\boldsymbol{\mu})$. The algorithm is summarized in Algorithm 3.

When the $K = |\mathcal{S}|$ is unknown, the MUSIC can also be applied by setting K to some integer larger than $|\mathcal{S}|$ in Algorithm 3. In such cases, the frequency set is determined by identifying the local maxima of $\mathcal{J}(\boldsymbol{\mu})$ using appropriate criteria to avoid numerical instabilities. In our experiments, we simply use the criterion that the amplitude of the local maxima $\hat{\mu}$ is larger than a preset threshold w > 0.

756	
757	Algorithm 3: multidimensional MUSIC
758	input : $y^{\epsilon}(t)$ sampled on $[-f, f]^{D}$ with $(N+1)^{D}$ sample points, K, L
759	1 Form the Hankel block matrix $H^{\epsilon} \in \mathbb{C}^{(L+1)^D \times (N-L+1)^D}$;
760	2 Perform the SVD: $H^{\epsilon} = \begin{bmatrix} U_1^{\epsilon} & U_2^{\epsilon} \end{bmatrix}$ diag $(\sigma_1^{\epsilon}, \cdots, \sigma_K^{\epsilon}, \cdots) \begin{bmatrix} V_1^{\epsilon} & V_2^{\epsilon} \end{bmatrix}^*$ where
761	$U_1 \in \mathbb{C}^{(L+1)^D imes K}$;
762	$\sigma = \sigma = 1$. Matrix $\sigma = \tau = \sigma = \tau$
763	3 Compute the MUSIC imaging function $\mathcal{J}(\mu)$ on the $\left[-\frac{4\pi}{2f}, \frac{4\pi}{2f}\right]$;
764	output: $S = \{K \text{ largest local maxima of } \mathcal{J}(\boldsymbol{\mu})\}$
166	

C COMPLEMENTS TO NUMERICAL RESULTS

C.1 CAPACITY OF LEARNING MODELS WITH LARGE MODEL ORDER

In this experiment, we perform two numerical tests to illustrate the capacity of learning mixture 771 model with a large model order in Algorithm 1 and 2. We first examine a 2-dimensional example of 772 a 12-component mixture model with a unified covariance matrix 0.3I. Using 1,000 samples from 773 this distribution, we compare the performance of Algorithm 1 and the EM algorithm in estimating 774 the component means. The EM algorithm is initialized with samples uniformly drawn from the 775 data and terminates when the log-likelihood increases less than 1×10^{-6} . The inputs of Algorithm 776 1 are set as $f_1 = f_2 = 3, L = 12, N = 25$. The results are shown in Figure 3. In the figure, the 777 true Gaussian components are illustrated as the red circles centered at the component mean with a 778 radius 1.5 times standard deviation, while the estimated ones are illustrated with the green dashed 779 circle. We observe that with the specific initialization used, the EM converges in 292 iterations but gets trapped in a local maxima. Algorithm 1 does not suffer from initialization issues and provides 781 a more accurate estimate of the mixture means.

782 Next, we consider a similar 12-component model but in a 100-dimensional space. The mixture means from Figure 3 are embedded into the \mathbb{R}^{100} and each mean is perturbed by a Gaussian vector 783 784 drawn from $\mathcal{N}(\mathbf{0}, 0.1 \mathbf{I}_{100})$. We apply Algorithm 2 and the EM algorithm to estimate the mix-785 ture means in this high-dimensional setting. The results are shown in Figure 4. For visualization 786 purposes, the estimates are projected onto the first two dimensions. It can also be seen from the es-787 timation error that the Algorithm 2 outperforms the EM algorithm under this setting. Furthermore, when considering only the 1-Wasserstein error in the first two dimensions, the Algorithm 2 shows 788 significantly better performance, with an error of 0.180 compared to the EM algorithm's error of 789 0.439.790

791 792

766 767

768 769

770

C.2 PROJECTION: ISSUES AND POTENTIAL SOLUTIONS

793 So far, we have focused on the numerical examples where the component means lie on or near a 794 2-dimensional subspace. However, a 2-dimensional projection may yield inaccurate estimations if the component means are distributed across a higher-dimensional space. The following experiment 796 illustrates this issue. In this experiment, we consider a 6-component mixture model in \mathbb{R}^3 . The 797 mixture means are $\{\pm Re_1, \pm Re_2, \pm Re_3\}$ where R = 4 and the covariance matrix is I_3 . We 798 draw 2,000 samples from this mixture model and use Algorithm 2 with k = 2 to estimate the 799 mixing distribution. The estimation results seem reasonable when examining the estimated means 800 projected onto the first two principal components. However, the accuracy degrades when considering the estimated means in the original \mathbb{R}^3 space. This discrepancy arises because, in this model, the 801 first two principal components span a subspace close to span $\{e_1, e_2\}$, making it challenging to 802 accurately estimate components whose means lie along the z-axis. As a result, projecting only onto 803 the subspace span $\{v_1, v_2\}$ makes it impossible to accurately estimate the third component. 804

To address this issue, one potential solution is to project the samples onto a higher-dimensional subspace and estimate the projected means. makes it impossible to accurately estimate the third component. As shown in (12), the time complexity of the *D*-dimensional MUSIC algorithm is exponential with respect to the data dimension *D*. Alternative multidimensional line spectral methods (e.g. Sarkar & Pereira (1995); Tang et al. (2014); Fei & Zhang (2023)) could also be applied, but they may encounter similar challenges. Another approach is to project the samples onto alternative



(a) Illustration of the EM algorithm with random initialization. Left: 1,000 samples(blue cross) drawn from a 12-component mixture model and the initialization means(black star) of the EM algorithm (converges in 292 iterations); Middle: the estimated components by the EM algorithm; Right: the estimated components by the EM algorithm (without samples illustrated).



(b) Illustration of the Algorithm 1. Left: imaging function values of the MUSIC algorithm and the 12 largest local maximal; Middle: the estimated components by the Algorithm 1; Right: the estimated components by the Algorithm 1 (without samples illustrated).





a 12-component mixture model and the initialization means(black star) of the EM algorithm (converges in 80 iterations); Middle: the estimated components by the EM algorithm; Right: the estimated components by the EM algorithm (without samples illustrated).



(b) Illustration of the Algorithm 2. Left: imaging function values of the MUSIC algorithm and the 12 largest local maximal; Middle: the estimated components by the Algorithm 1; Right: the estimated components by the Algorithm 1 (without samples illustrated).





D.1 COMPUTATIONAL RESOLUTION LIMIT

Consider the Fourier measurements of the high-dimensional line spectral signal as (13) and assume that $\|\epsilon(t)\|_{\infty} < \sigma$. The following theorem gives an upper bound for the computational resolution limit for the number detection:

Theorem 4. (*Liu & Zhang (2021a), Theorem 2.3*) Let the Fourier measurement (13) be generated by an n-sparse measure $\nu = \sum_{i=1}^{K} w_i \delta_{\mu_i}, \mu_i \in B^D_{\frac{(K-1)w}{2f}}(\mathbf{0})$. Let $K \ge 2$ and assume the following separation condition is satisfied

$$\Delta = \min_{1 \le i < j \le K} \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\| \ge \frac{C_2(K, D)}{f} \left(\frac{\sigma}{w_{\min}}\right)^{\frac{1}{2K-2}}$$
(30)

where $C_2(K,D) = 4.4we(w/2)^{s-1}(K(K-1)/w)^{\xi(s-1)}$ with

$$\xi(k) = \left\{ \begin{array}{ll} \sum_{i=1}^{k} \frac{1}{i}, & k \ge 1 \\ 0, & k = 0, \end{array} \right.$$

and s being the the dimension of the smallest subspace in \mathbb{R}^D which contains the set $\{\mu_1, \dots, \mu_K\}$. Then there do not exist any σ -admissible measures of y(t) with less than n components. This theorem provides an upper bound of $\mathcal{R}_{K,D}$. The lower bound has also been characterized by the following proposition:

Proposition 2. (Liu & Zhang (2021a), Proposition 2.4) For given $0 < \sigma < w_{\min}$ and $K \ge 2$, there exist an K-sparse measure in \mathbb{R}^D , $\nu = \sum_{i=1}^K w_i \delta_{\mu_i}$ and an (n-1)-sparse measure in \mathbb{R}^D , $\hat{\nu} = \sum_{i=1}^{K-1} \hat{w}_i \delta_{\hat{\mu}_i}$, such that $\|\mathcal{F}\hat{\nu}(t) - \mathcal{F}\nu(t)\|_{\infty} < \sigma$, $\|t\|_2 \le f$. Moreover

$$\min_{1 \le i \le K} |w_i| = w_{\min}, \quad \min_{1 \le i < j \le K} \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_2 = \frac{C_1(K, D)}{f} \left(\frac{\sigma}{w_{\min}}\right)^{\frac{1}{2K-2}}.$$

where $C_1(K, D) = 0.81e^{-\frac{3}{2}}$.

The above results indicate that

$$\frac{C_1(K,D)}{f} \left(\frac{\sigma}{w_{\min}}\right)^{\frac{1}{2K-2}} < \mathcal{R}_{K,D} \le \frac{C_2(K,D)}{f} \left(\frac{\sigma}{w_{\min}}\right)^{\frac{1}{2K-2}}.$$
(31)

The computational resolution limit for the support recovery of (13) has also been established in Liu & Zhang (2021a). Denote the computational resolution limit for support recovery as $\tilde{\mathcal{R}}_{K,D}$ and the results indicate that

$$\frac{\tilde{C}_1(K,D)}{f} \left(\frac{\sigma}{w_{\min}}\right)^{\frac{1}{2K-1}} < \tilde{\mathcal{R}}_{K,D} \le \frac{\tilde{C}_2(K,D)}{f} \left(\frac{\sigma}{w_{\min}}\right)^{\frac{1}{2K-1}},\tag{32}$$

where $\tilde{C}_1(K,D) = 0.49e^{-\frac{3}{2}}$ and $\tilde{C}_2(K,D) = 5.88\pi e 4^{D-1} \left((K+2)(K-1)/2 \right)^{\xi(D-1)}$. From (31) and (32), it reveals the difference between these two tasks quantitatively by the $\frac{1}{2K-2}$ and $\frac{1}{2K-1}$ powered on the signal noise ratio term σ/w_{\min} .

1000 D.2 PROOF OF THEOREM 1

1002 *Proof.* By setting $\epsilon = w_{\min} \left(\frac{\Delta f}{C_2(K,D)}\right)^{2K-2}$ in Proposition 1, we see that for

$$n \ge C_{K,D} \log\left(\frac{4}{\delta}\right) \frac{e^{2f^2 \sigma_{\max}(\boldsymbol{\Sigma})}}{w_{\min}^2 (f\Delta)^{4K-4}}$$

where $C_{K,D} = 4 (C_2(K,D))^{4K-4}$, we have

 $\|\epsilon_n(\boldsymbol{t})\|_{\infty} \le w_{\min}\left(\frac{\Delta f}{C_2(K,D)}\right)^{2K-2}, \|\boldsymbol{t}\|_2 \le f$

1012 with probability at least $1 - \delta$. The rest follows from Theorem (4).

1016 Similar to Theorem 1, the sample size requirement for estimating the means is given by

Theorem 5. Consider the D-dimensional mixture model $\sum_{i=1}^{K} w_i \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ with $\boldsymbol{\mu}_i \in B_{\frac{(K-1)w}{2f}}^D(\mathbf{0})$. For any $\delta \in (0, 1)$, if the sample size n satisfies that

$$n \ge \tilde{C}_{K,D} \log\left(\frac{4}{\delta}\right) \frac{e^{2f^2 \sigma_{\max}(\mathbf{\Sigma})}}{w_{\min}^2 (f\Delta)^{4K-2}}.$$
(33)

1024 Then with probability $1 - \delta$, $\Delta \ge \tilde{\mathcal{R}}_{D,K}$ holds. Here $C_{K,D}$ is a constant only relying on K and D.

The proof of the theorem is the same as that of Theorem 1.

1026 E PROOF OF PROPOSITION 1

Proof. Note that

$$e^{t^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{t}} \psi_n(\boldsymbol{t}) = \frac{1}{n} \sum_{j=1}^n e^{t^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{t}} \cos\langle \boldsymbol{x}_j, \boldsymbol{t} \rangle + \iota \frac{1}{n} \sum_{j=1}^n e^{t^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{t}} \sin\langle \boldsymbol{x}_j, \boldsymbol{t} \rangle.$$

Applying the Hoeffding's inequality to the real and imaginary parts, we have

$$\mathbb{P}\left(\left|\frac{1}{n}\sum_{j=1}^{n}e^{t^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{t}}\cos\langle\boldsymbol{x}_{j},\boldsymbol{t}\rangle-e^{t^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{t}}\mathfrak{R}\left(\boldsymbol{\phi}(\boldsymbol{t})\right)\right|>\epsilon\right)\leq2\exp\left(-\frac{n\epsilon^{2}}{2e^{2t^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{t}}}\right),\\\mathbb{P}\left(\left|\frac{1}{n}\sum_{j=1}^{n}e^{t^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{t}}\sin\langle\boldsymbol{x}_{j},\boldsymbol{t}\rangle-e^{t^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{t}}\Im\left(\boldsymbol{\phi}(\boldsymbol{t})\right)\right|>\epsilon\right)\leq2\exp\left(-\frac{n\epsilon^{2}}{2e^{2t^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{t}}}\right).$$

Hence,

$$\begin{split} & \mathbb{P}\left(\left|e^{t^{\mathrm{T}}\Sigma t}\psi_{n}(t) - \sum_{i=1}^{K}w_{i}\exp\left(\iota\langle\boldsymbol{\mu}_{i}, t\rangle\right)\right| > \epsilon\right) = \mathbb{P}\left(\left|e^{t^{\mathrm{T}}\Sigma t}[\psi_{n}(t) - \phi(t)]\right| > \epsilon\right) \\ & \leq \mathbb{P}\left(\left|\Re\left(e^{t^{\mathrm{T}}\Sigma t}[\psi_{n}(t) - \phi(t)]\right)\right| > \frac{\epsilon}{\sqrt{2}}\right) + \mathbb{P}\left(\left|\Im\left(e^{t^{\mathrm{T}}\Sigma t}[\psi_{n}(t) - \phi(t)]\right)\right| > \frac{\epsilon}{\sqrt{2}}\right) \\ & = \mathbb{P}\left(\left|\frac{1}{n}\sum_{j=1}^{n}e^{t^{\mathrm{T}}\Sigma t}\cos\langle x_{j}, t\rangle - e^{t^{\mathrm{T}}\Sigma t}\Re\left(\phi(t)\right)\right| > \frac{\epsilon}{\sqrt{2}}\right) + \mathbb{P}\left(\left|\frac{1}{n}\sum_{j=1}^{n}e^{t^{\mathrm{T}}\Sigma t}\sin\langle x_{j}, t\rangle - e^{t^{\mathrm{T}}\Sigma t}\Im\left(\phi(t)\right)\right| > \frac{\epsilon}{\sqrt{2}}\right) \\ & \leq 4\exp\left(-\frac{n\epsilon^{2}}{4e^{2t^{\mathrm{T}}\Sigma t}}\right) \leq 4\exp\left(-\frac{n\epsilon^{2}}{4e^{2\|t\|_{2}^{2}\sigma_{\max}(\Sigma)}}\right) < \delta. \end{split}$$
where we used $n > 4\log\left(\frac{4}{\delta}\right) \frac{e^{2\|t\|_{2}^{2}\sigma_{\max}(\Sigma)}}{\epsilon^{2}}$ in the last inequality. \Box

1055 E.1 PROOF OF THEOREM 2

Proof. Notice that

$$\left\|\operatorname{Proj}_{\boldsymbol{W}_{k}}\boldsymbol{X}\right\|_{2}^{2} = \sum_{j=1}^{n} \left\|\operatorname{Proj}_{\boldsymbol{W}_{k}}\boldsymbol{x}_{j}\right\|_{2}^{2} = \sum_{j=1}^{n} \sum_{l=1}^{k} |\langle \boldsymbol{x}_{j}, \boldsymbol{w}_{l} \rangle|^{2}.$$

Taking the expectation, we get

$$\begin{split} \mathbb{E} \left\| \operatorname{Proj}_{\boldsymbol{W}_{k}} \boldsymbol{X} \right\|_{2}^{2} &= \sum_{j=1}^{n} \sum_{l=1}^{k} \mathbb{E} |\langle \boldsymbol{x}_{j}, \boldsymbol{w}_{l} \rangle|^{2} \\ &= \sum_{j=1}^{n} \sum_{l=1}^{k} \sum_{i=1}^{K} \mathbb{E} \left[|\langle \boldsymbol{x}_{j}, \boldsymbol{w}_{l} \rangle|^{2} | \boldsymbol{z}_{j} = i \right] \mathbb{P}(\boldsymbol{z}_{j} = i) \\ &= \sum_{j=1}^{n} \sum_{l=1}^{k} \sum_{i=1}^{K} w_{i} \left(\sigma^{2} + \mathbb{E} \left[\langle \boldsymbol{x}_{j}, \boldsymbol{w}_{l} \rangle | \boldsymbol{z}_{j} = i \right]^{2} \right) \\ &= \sum_{j=1}^{n} \sum_{l=1}^{k} \sum_{i=1}^{K} w_{i} \left(\sigma^{2} + \langle \boldsymbol{\mu}_{i}, \boldsymbol{w}_{l} \rangle^{2} \right) \\ &= \sum_{j=1}^{n} \left(k\sigma^{2} + \sum_{i=1}^{K} w_{i} \sum_{l=1}^{k} \langle \boldsymbol{\mu}_{i}, \boldsymbol{w}_{l} \rangle^{2} \right) \\ &= n \left(k\sigma^{2} + \sum_{i=1}^{K} w_{i} \left\| \operatorname{Proj}_{\boldsymbol{W}_{k}} \boldsymbol{\mu}_{i} \right\|_{2}^{2} \right) \end{split}$$

Case 1: k = K. We have

$$\mathbb{E}\left\|\operatorname{Proj}_{\boldsymbol{W}_{K}}\boldsymbol{X}\right\|_{2}^{2} \leq n\left(k\sigma^{2} + \sum_{i=1}^{K} w_{i} \left\|\boldsymbol{\mu}_{i}\right\|_{2}^{2}\right) = \mathbb{E}\left\|\operatorname{Proj}_{\boldsymbol{U}_{K}}\boldsymbol{X}\right\|_{2}^{2}$$

1084 where $U_K = \text{span}\{\mu_1, \cdots \mu_K\}.$

1086 **Case 2:** k < K. We show that the k-dimensional subspace W_k maximizing the $\mathbb{E} \| \operatorname{Proj}_{W_k} X \|_2$ is 1087 the subspace of U_K . Notice that

1088 1089 1090

1091

1095 1096

1098 1099

1100 1101 1102

1103 1104

1105

1109 1110 1111

1115

$$\sum_{i=1}^{K} w_i \left\| \operatorname{Proj}_{\boldsymbol{W}_k} \boldsymbol{\mu}_i \right\|_2^2 = \sum_{i=1}^{K} \left\| \operatorname{Proj}_{\boldsymbol{W}_k} \sqrt{w_i} \boldsymbol{\mu}_i \right\|_2^2$$
$$= \left\| \operatorname{Proj}_{\boldsymbol{W}_k} \left[\sqrt{w_1} \boldsymbol{\mu}_1 \quad \cdots \quad \sqrt{w_K} \boldsymbol{\mu}_K \right] \right\|_2^2.$$

¹⁰⁹³ Therefore, the k-dimensional subspace maximizing the projection above is the subspace spanned by the first k right eigenvectors of the SVD of $[\sqrt{w_1}\mu_1 \cdots \sqrt{w_K}\mu_K]$. This subspace W_k satisfies

$$\boldsymbol{W}_k \subseteq Im([\sqrt{w_1}\boldsymbol{\mu}_1 \quad \cdots \quad \sqrt{w_K}\boldsymbol{\mu}_K]) = \boldsymbol{U}_K.$$

Case 3: k > K. We prove that the k-dimensional subspace W_k maximizing $\mathbb{E} \| \operatorname{Proj}_{W_k} X \|_2$ must contain U_K . Indeed, for any W_k such that $U_K \subset W_k$, we have

$$\mathbb{E} \left\| \operatorname{Proj}_{\boldsymbol{W}_{k}} \boldsymbol{X} \right\|_{2} = n \left(k \sigma^{2} + \sum_{i=1}^{K} w_{i} \left\| \boldsymbol{\mu}_{i} \right\|_{2}^{2} \right).$$

F QUADRATIC PROGRAMMING OPTIMIZATION

In this section, we introduce the quadratic programming(QP) optimization applied in Algorithm 1 for recovering the component weights. The general formulation of the QP can be expressed as

minimize
$$\frac{1}{2} \boldsymbol{x}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{x} + \boldsymbol{q}^{\mathrm{T}} \boldsymbol{x} + r$$

subject to $\boldsymbol{G} \boldsymbol{x} \leq \boldsymbol{h}, \quad \boldsymbol{A} \boldsymbol{x} = \boldsymbol{b}$ (34)

where $P \in \mathbb{R}^{n \times n}$, $G \in \mathbb{R}^{m \times n}$, $A \in \mathbb{R}^{p \times n}$ and P is positive-definite. This optimization program can be viewed as minimizing a convex quadratic function over a polyhedron. For a more comprehensive introduction of the QP optimization, we refer to Boyd & Vandenberghe (2004).

1116 F.1 MIXING WEIGHTS ESTIMATION

After achieving the model order \hat{K} and mean set $\{\hat{\mu}_i : 1 \le i \le \hat{K}\}$, the corresponding weights are estimated by solving

1120
1121
1122
1123
1124
1125
minimize
$$\left\| e^{-t^{\mathrm{T}} \Sigma t} \sum_{i=1}^{\hat{K}} w_i e^{i \langle \hat{\mu}_i, t \rangle} - \psi_n(t) \right\|_2$$

subject to $w_i \ge w$, $\sum_{i=1}^{\hat{K}} w_i = 1$
(35)

1126 The program (35) can be reformulated as a quadratic programming(QP) optimization and can be 1127 efficiently solved by well-established convex optimization toolboxes². Next, we show how to fit the 1128 optimization problem (35) into the framework of (34). To simplify the notation, we replace \hat{K} , $\hat{\mu}_i$'s 1129 with the unhatted ones. Notice that we can write

²In the numerical experiments, we use the python package cvxpy to implement the quadratic programming.

where $\boldsymbol{A} \in \mathbb{C}^{(N+1)^2 \times K}$ and $\boldsymbol{b} \in \mathbb{C}^{(N+1)^2}$ such that

 $\boldsymbol{A}_{j,i} = e^{\iota \langle \boldsymbol{\mu}_i, \boldsymbol{t}_j \rangle}, \quad \boldsymbol{b}_j = e^{\boldsymbol{t}_j^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{t}_j} \psi_n(\boldsymbol{t}_j).$

~

Here t_j is the j-th component of the vectorized sample points. Then the objective function can be further written as

$$\begin{aligned} & \prod_{i=1}^{1140} \left\| \sum_{i=1}^{K} w_{i} e^{\iota \langle \boldsymbol{\mu}_{i}, \boldsymbol{t} \rangle} - e^{\boldsymbol{t}^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{t}} \psi_{n}(\boldsymbol{t}) \right\|_{2}^{2} = \|\Re(\boldsymbol{A}\pi - \boldsymbol{b})\|_{2}^{2} + \|\Im(\boldsymbol{A}\pi - \boldsymbol{b})\|_{2}^{2} \\ & \prod_{i=1}^{1142} \left\| \Re(\boldsymbol{A}) \pi - \Re(\boldsymbol{b}) \|_{2}^{2} + \|\Im(\boldsymbol{A}) \pi - \Im(\boldsymbol{b})\|_{2}^{2} \\ & = \|\Re(\boldsymbol{A}) \pi - \Re(\boldsymbol{b})\|_{2}^{2} + \|\Im(\boldsymbol{A}) \pi - \Im(\boldsymbol{b})\|_{2}^{2} \\ & = (\Re(\boldsymbol{A}) \pi - \Re(\boldsymbol{b}))^{\mathrm{T}}(\Re(\boldsymbol{A}) \pi - \Re(\boldsymbol{b})) + (\Im(\boldsymbol{A}) \pi - \Im(\boldsymbol{b}))^{\mathrm{T}}(\Im(\boldsymbol{A}) \pi - \Im(\boldsymbol{b})) \\ & = \pi^{\mathrm{T}}[\Re(\boldsymbol{A})^{\mathrm{T}} \Re(\boldsymbol{A}) + \Im(\boldsymbol{A})^{\mathrm{T}} \Im(\boldsymbol{A})] \pi - 2[\Re(\boldsymbol{b})^{\mathrm{T}} \Re(\boldsymbol{A}) + \Im(\boldsymbol{b})^{\mathrm{T}} \Im(\boldsymbol{A})] \pi + \Re(\boldsymbol{b})^{\mathrm{T}} \Re(\boldsymbol{b}) + \Im(\boldsymbol{b})^{\mathrm{T}} \Im(\boldsymbol{b}) . \end{aligned}$$

Therefore, we can fit (35) into the QP framework by setting

1149
1149

$$P = \Re (A)^{\mathrm{T}} \Re (A) + \Im (A)^{\mathrm{T}} \Im (A), \quad q = \Re (A)^{\mathrm{T}} \Re (b) + \Im (A)^{\mathrm{T}} \Im (b), \quad r = \frac{1}{2} \Re (b)^{\mathrm{T}} \Re (b) + \frac{1}{2} \Im (b)^{\mathrm{T}} \Im (b)$$
1150

in the objective function and setting

$$G = -I_K$$
, $h = -w \mathbf{1}_{K \times 1}$, $A = \mathbf{1}_{1 \times K}$, $b = 1$.

G THE EM ALGORITHM

In this section, we describe the EM algorithm used in the numerical experiments for comparison with our algorithms.

	Algorithm	5: The	EM	algorithm	(Fixed	Covariance	Matrix)
--	-----------	--------	----	-----------	--------	------------	---------

input : samples x_1, \dots, x_n , model order k, covariance matrix Σ , initial guess \hat{w}_i 's, $\hat{\mu}_i$'s

1 *Expectation Step:* For $i = 1, \dots, k$, compute

$$\gamma_i^j = rac{w_i g(oldsymbol{x}_i; \hat{oldsymbol{\mu}}_i, oldsymbol{\Sigma})}{\sum_{i=1}^k g(oldsymbol{x}_i; \hat{oldsymbol{\mu}}_i, oldsymbol{\Sigma})}, \quad j = 1, \cdots, n.$$

2 Maximization Step: Compute the weights and weighted means:

$$\hat{w}_i = rac{1}{n} \sum_{j=1}^n \gamma_i^j, \quad \hat{\mu}_i = rac{\sum_{j=1}^n \gamma_i^j \boldsymbol{x}_j}{\sum_{j=1}^n \gamma_i^j}, \quad i = 1, \cdots, k.$$

3 Iterate steps 1 and 2 until convergence.

In the numerical tests, we assume that the covariance matrix Σ is known as prior information. If the covariance matrix is unknown, it is updated in the maximization step by

$$\hat{\boldsymbol{\Sigma}} = \frac{\sum_{i=1}^{k} \sum_{j=1}^{n} \gamma_i^j (\boldsymbol{x}_j - \hat{\boldsymbol{\mu}}_i) (\boldsymbol{x}_j - \hat{\boldsymbol{\mu}}_i)^{\mathrm{T}}}{n}$$

for the unified covariance matrix case and

$$\hat{\boldsymbol{\Sigma}}_i = rac{\sum_{j=1}^n \gamma_i^j (\boldsymbol{x}_j - \hat{\boldsymbol{\mu}}_i) (\boldsymbol{x}_j - \hat{\boldsymbol{\mu}}_i)^{\mathrm{T}}}{\sum_{i=1}^n \gamma_i^j}$$

for the general Gaussian mixture model.