A Gaussian Mixture Model Representation of Endmember Variability in Hyperspectral Unmixing

Yuan Zhou[®], Student Member, IEEE, Anand Rangarajan, Member, IEEE, and Paul D. Gader, Fellow, IEEE

Abstract—Hyperspectral unmixing while considering endmember variability is usually performed by the normal compositional model, where the endmembers for each pixel are assumed to be sampled from unimodal Gaussian distributions. However, in real applications, the distribution of a material is often not Gaussian. In this paper, we use Gaussian mixture models (GMM) to represent endmember variability. We show, given the GMM starting premise, that the distribution of the mixed pixel (under the linear mixing model) is also a GMM (and this is shown from two perspectives). The first perspective originates from random variable transformations and gives a conditional density function of the pixels given the abundances and GMM parameters. With proper smoothness and sparsity prior constraints on the abundances, the conditional density function leads to a standard maximum a posteriori (MAP) problem which can be solved using generalized expectation maximization. The second perspective originates from marginalizing over the endmembers in the GMM, which provides us with a foundation to solve for the endmembers at each pixel. Hence, compared to the other distribution based methods, our model can not only estimate the abundances and distribution parameters, but also the distinct endmember set for each pixel. We tested the proposed GMM on several synthetic and real datasets, and showed its potential by comparing it to current popular methods.

12

13

14

16

17

18

21

22

27

28

29

30

31

32

Index Terms—Endmember extraction, endmember variability, hyperspectral image analysis, linear unmixing, Gaussian mixture model.

I. INTRODUCTION

THE formation of hyperspectral images can be simplified by the *linear mixing model* (LMM), which assumes that the physical region corresponding to a pixel contains several pure materials, so that each material contributes a fraction of its spectra based on area to the final spectra of the pixel. Hence, the observed spectra $\mathbf{y}_n \in \mathbb{R}^B$, $n = 1, \ldots, N$ (B is the number of wavelengths and N is the number of pixels) is a (non-negative) linear combination of the pure material

Manuscript received June 30, 2017; revised November 23, 2017; accepted January 10, 2018. This work was supported by NSF IIS under Grant 1743050. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Jocelyn Chanussot. (Corresponding author: Yuan Zhou.)

The authors are with the Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL 32611 USA (e-mail: yuan@cise.ufl.edu; anand@cise.ufl.edu; pgader@cise.ufl.edu).

This paper has supplementary downloadable material available a http://ieeexplore.ieee.org., provided by the author.

Color versions of one or more of the figures in this paper are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/TIP.2018.2795744

(called *endmember*) spectra $\mathbf{m}_j \in \mathbb{R}^B$, j = 1, ..., M (M is the number of endmembers), i.e.

$$\mathbf{y}_{n} = \sum_{j=1}^{M} \mathbf{m}_{j} \alpha_{nj} + \mathbf{n}_{n}, \text{ s.t. } \alpha_{nj} \ge 0, \quad \sum_{j=1}^{M} \alpha_{nj} = 1, \quad (1)$$

where a_{nj} is the proportion (called *abundance*) for the *j*th endmember at the *n*th pixel (with the positivity and sumto-one constraint) and $\mathbf{n}_n \in \mathbb{R}^B$ is additive noise. Here, the endmember set $\{\mathbf{m}_j : j=1,\ldots,M\}$ is fixed for all the pixels. This model simplifies the unmixing problem to a matrix factorization one, leading to efficient computation and simple algorithms such as iterative constrained endmembers (ICE), vertex component analysis (VCA), piecewise convex multiplemodel endmember detection (PCOMMEND) [1]–[3] etc., which receive comprehensive reviews in [4] and [5].

46

50

55

61

65

73

74

76

However, in practice the LMM may not be valid in many real scenarios. Even for a *pure* pixel that only contains one material, its spectrum may not be consistent over the whole image. This is due to several factors such as atmospheric conditions, topography and intrinsic variability. For example, in vegetation, multiple scattering and biotic variation (e.g. differences in biochemistry and water content) cause different reflectances among the same species. For urban scenes, the incidence and emergence angles could be different for the same roof, causing different reflectances. For minerals, the spectroscopy model developed by Hapke also considers the porosity and roughness of the material as variable [6].

In the first and third example above, Eq. (1) can be generalized to a more abstract form $\mathbf{y}_n = F(\{\mathbf{m}_j, \alpha_{nj} : j=1,\ldots M\})$, which leads to nonlinear mixing models. For example, Halimi et al. [7] used bilinear models to handle the vegetation case, which was also investigated using several different nonlinear functions [8]. In [9], the Hapke model was used to model intimate interaction among minerals. There are also works that use kernels for flexible nonlinear mixing [10], [11]. A panoply of nonlinear models can be found in the review article [12]. We note that in these models, a fixed endmember set is still assumed while using a more complicated unmixing model.

While nonlinear models abound lately, it is still difficult to account for all the scenarios. On the contrary, the LMM still has physical significance with the intuitive area assumption. To model real scenarios more accurately, researchers have

87

91

93

97

100

101

102

103

104

106

108

110

112

114

115

116

117

119

120

121

123

125

127

129

130

131

taken another route by generalizing Eq. (1) to

$$\mathbf{y}_n = \sum_{i=1}^M \mathbf{m}_{nj} \alpha_{nj} + \mathbf{n}_n, \tag{2}$$

where $\{\mathbf{m}_{nj} \in \mathbb{R}^B : j = 1, ..., M\}, n = 1, ..., N$ could be different for each n, i.e. the endmember spectra for each pixel could be different. This is called endmember variability, and has also received a lot of attention in the community [13], [14]. Note that given $\{\mathbf{y}_n\}$, inferring $\{\mathbf{m}_{nj}, \alpha_{nj}\}$ is a much more difficult problem than inferring $\{\mathbf{m}_j, \alpha_{nj}\}$ in Eq. (1). Hence, in many papers $\{\mathbf{m}_{nj}\}$ are assumed to be from a spectral library, which is usually called *supervised unmixing* [15]–[17]. On the other hand, if the endmember spectra are to be extracted from the image, we call them unsupervised unmixing models [18]–[20]. Obviously, unsupervised unmixing more challenging than its supervised counterpart and hence more assumptions are used in this case, such as the spatial smoothness of abundances and endmember variability [21]-[23], small mutual distance between the endmembers [22], small magnitude or spectral smoothness of the endmember variability [22], [23].

We can also categorize the papers on endmember variability by how this variability is modeled. In the review paper [14], it can be modeled as a endmember set [17], [20] or as a distribution [24]-[26]. One of the widely used set based methods is multiple endmember spectral mixture analysis (MESMA) [17], which tries every endmember combination and selects the one with the smallest error. There are many variations to the original MESMA. For example, the multiple-endmember linear spectral unmixing model (MELSUM) solves the linear equations directly using the pseudo-inverse and discards the solutions with negative abundances [27]; automatic Monte Carlo unmixing (AutoMCU) picks random combinations for unmixing and averages the resulting abundances as the final results [28], [29]. Besides MESMA variants, there are also many other set based methods. For example, endmember bundles form bundles from automated extracted endmembers, take minimum and maximum abundances from bundle based unmixing, and average them as final abundances [20]; sparse unmixing imposes a sparsity constraint on the abundances based on endmembers composed of all spectra from the spectral library [30]. A comprehensive review can be found in [13] and [14]. One disadvantage of set based methods is that their complexity increases exponentially with increasing library size hence in practice a laborious library reduction approach may be required [31].

The distribution based approaches assume that the endmembers for each pixel are sampled from probability distributions [e.g. Gaussian, a.k.a. normal compositional model (NCM)], and hence embrace large libraries while being numerically tractable [15], [32]. Here, we give an overview of NCM because of its simplicity and popularity [16], [18], [19]. Suppose the jth endmember at the nth pixel follows a Gaussian distribution $p\left(\mathbf{m}_{nj}\right) = \mathcal{N}\left(\mathbf{m}_{nj}|\boldsymbol{\mu}_{j},\boldsymbol{\Sigma}_{j}\right)$ where $\boldsymbol{\mu}_{j} \in \mathbb{R}^{B}$ and $\boldsymbol{\Sigma}_{j} \in \mathbb{R}^{B \times B}$, and the additive noise also follows a Gaussian distribution $p\left(\mathbf{n}_{n}\right) = \mathcal{N}\left(\mathbf{n}_{n}|\mathbf{0},\mathbf{D}\right)$ where \mathbf{D} is the noise covariance matrix. The random variable transformation (r.v.t.)

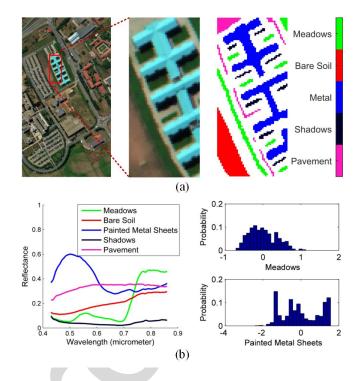


Fig. 1. (a) Original Pavia University image and selected ROI with its ground truth image. (b) Mean spectra of the identified 5 endmembers and histograms of meadows and painted metal sheets (shadow is termed as endmember to conform with the LMM though the area under shadow can be any material). PCA is used to project the multidimensional pixels to single values which are counted in the histograms. Although the histogram of meadows may appear to be a Gaussian distribution, that of painted metal sheets is obviously neither a unimodal Gaussian or Beta distribution.

(2) suggests that the probability density function of \mathbf{y}_n can be derived as

$$p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right) = \mathcal{N}\left(\mathbf{y}_{n}|\sum_{j=1}^{M}\alpha_{nj}\boldsymbol{\mu}_{j},\sum_{j=1}^{M}\alpha_{nj}^{2}\boldsymbol{\Sigma}_{j} + \mathbf{D}\right), \quad (3)$$

137

139

141

143

144

145

146

147

148

149

150

151

152

154

where $\alpha_n := [\alpha_{n1}, \dots, \alpha_{nM}]^T$, $\Theta := \{\mu_j, \Sigma_j : j = 1, \dots, M\}$. The conditional density function in (3) is usually embedded in a Bayesian framework such that we can incorporate priors and also estimate hyperparameters. Then, NCM uses different optimization approaches, e.g. expectation maximization [32], sampling methods [18], [19], [25], particle swarm optimization [24], to determine the parameters $\{\mu_i, \Sigma_i\}$ and $\{\alpha_{ni}\}$.

There are few papers that use other distributions. In [15], X. Du *et al.* note that the Gaussian distribution may allow negative values which are not realistic. In addition, the real distribution may be skewed. Hence, they introduce a Beta compositional model (BCM) to model the variability. The problem is that the true distribution may not be well approximated by any unimodal distribution. Consider the Pavia University dataset shown in Fig. 1, where the multidimensional pixels are projected to one dimension to afford better visualization. Among the manually identified materials, we can see that although the histogram of meadows may look like a Gaussian distribution, that of painted metal sheets has multiple peaks and cannot be approximated by either a Gaussian or

204

205

207

211

212

213

215

221

223

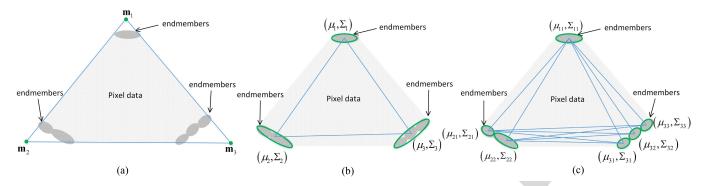


Fig. 2. Comparison of the mechanisms among LMM, NCM and GMM. We have 3 endmembers represented by the darken gray areas. LMM tries to find a set of endmembers that fit the pixel data. NCM tries to find a set of Gaussian centers that fit the pixel data, with error weighted by the covariance matrices. GMM tries to find Gaussian centers such that all their linear combinations fit the pixel data, with each weighted by the prior π_k . We may use 6 endmembers with NCM, but then the prior information is lost.

Beta distribution. This is due to different angles of these sheets on the roof. Since each piece of metal sheet is tilted, it forms a cluster of reflectances which contributes to a peak in the histogram. This example shows that we should use a more flexible distribution to represent the endmember variability.

158

160

161

162

163

164

165

167

169

171

173

174

175

178

180

181

182

184

185

186

188

189

190

191

192

193

194

195

197

AO:1

In this paper, we use a mixture of Gaussians to approximate any distribution that an endmember may exhibit, and solve the LMM by considering endmember variability. In a nutshell, the Gaussian mixture model (GMM) models $p(\mathbf{m}_{ni})$ by a mixture of Gaussians, say $p(\mathbf{m}_{ni}) =$ $\sum_{k} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right)$, and then obtains the distribution of y_n by the r.v.t. (2), which turns out to be another mixture of Gaussians and can be used for inference of the unknown parameters. Here, we briefly explain how GMM works intuitively by comparing it to the NCM with the details given later. The maximum likelihood estimate (MLE) of NCM (using (3)) aims to find $\{\mu_i\}$ such that its linear combination matches \mathbf{y}_n . Contrary to NCM, GMM aims to find $\{\mu_{ik}\}$ such that all of its linear combinations match y_n . Suppose we have μ_{11} , μ_{21} , μ_{22} , μ_{31} , μ_{32} , μ_{33} : then there are 6 combinations as explained in Fig. 2, but with emphasis weighted by $\{\pi_{jk}\}$ which determines the prior probability of each linear combination.

Based on the GMM formulation, we propose a supervised version and an unsupervised version for unmixing. The supervised version takes a library as input and estimates the abundances. The unsupervised version assumes that there are regions of pure pixels, hence segments the image first to get pure pixels and then performs unmixing. Another advantage over the other distribution based methods is that we can also estimate the endmembers for each pixel, which is not achievable by NCM or BCM. Note that estimating endmembers for each pixel is generally common in non-distribution methods, both from the signal processing community [21]–[23] or the remote sensing community [17], [27]. But it is often achieved in the context of least-squares based unmixing [33]–[35], unlike what we propose here using distribution based unmixing.

Notation: As usual, $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the multivariate Gaussian density function with center $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a matrix with m rows and n columns. The Hadamard product of two matrices (elementwise multiplication) is denoted by \circ while the Kronecker

product is denoted by \otimes . (A) $_{jk}$ denotes the element at the jth row and kth column of matrix \mathbf{A} . (A) $_j$ denotes the jth row of \mathbf{A} transposed (treating \mathbf{A} as a vector), i.e. for $\mathbf{A} = [\mathbf{a}_1, \dots \mathbf{a}_n]^T$, (A) $_j = \mathbf{a}_j$. vec (A) denotes the vectorization of \mathbf{A} , i.e. concatenating the columns of \mathbf{A} . $\delta_{jk} = 1$ when j = k and 0 otherwise. $\mathbb{E}_{\mathbf{x}}(f(\mathbf{x}))$ is the expected value of $f(\mathbf{x})$ given random variable \mathbf{x} . We use $i = \sqrt{-1}$ instead of as an index throughout the paper.

II. MATHEMATICAL PRELIMINARIES

A. Linear Combination of GMM Random Variables

To use the Gaussian mixture model to model endmember variability, we start by assuming that \mathbf{m}_{nj} follows a Gaussian mixture model (GMM) and the noise also follows a Gaussian distribution. The distribution of \mathbf{y}_n is obtained using the following theorem.

Theorem 1: If the random variable \mathbf{m}_{nj} has a density function

$$p\left(\mathbf{m}_{nj}|\boldsymbol{\Theta}\right) := f_{\mathbf{m}_{j}}\left(\mathbf{m}_{nj}\right) = \sum_{k=1}^{K_{j}} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj}|\boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right), \quad (4) \quad {}_{210}$$

s.t. $\pi_{jk} \geq 0$, $\sum_{k=1}^{K_j} \pi_{jk} = 1$, with K_j being the number of components, π_{jk} ($\mu_{jk} \in \mathbb{R}^B$ or $\Sigma_{jk} \in \mathbb{R}^{B \times B}$) being the weight (mean or covariance matrix) of its kth Gaussian component, $\Theta := \{\pi_{jk}, \mu_{jk}, \Sigma_{jk} : j = 1, \dots, M, k = 1, \dots, K_j\}$, $\{\mathbf{m}_{nj} : j = 1, \dots, M\}$ are independent, and the random variable \mathbf{n}_n has a density function $p(\mathbf{n}_n) := \mathcal{N}(\mathbf{n}_n | \mathbf{0}, \mathbf{D})$, then the density function of \mathbf{y}_n given by the r.v.t. $\mathbf{y}_n = \sum_{j=1}^M \mathbf{m}_{nj} \alpha_{nj} + \mathbf{n}_n$ is another GMM

$$p(\mathbf{y}_n|\boldsymbol{\alpha}_n,\boldsymbol{\Theta},\mathbf{D}) = \sum_{\mathbf{k}\in\mathcal{K}} \pi_{\mathbf{k}} \mathcal{N}(\mathbf{y}_n|\boldsymbol{\mu}_{n\mathbf{k}},\boldsymbol{\Sigma}_{n\mathbf{k}}),$$
 (5)

where $K := \{1, ..., K_1\} \times \{1, ..., K_2\} \times \cdots \times \{1, ..., K_M\}$ is the Cartesian product of the M index sets, $\mathbf{k} := (k_1, ..., k_M) \in \mathcal{K}$, $\pi_{\mathbf{k}} \in \mathbb{R}$, $\mu_{n\mathbf{k}} \in \mathbb{R}^B$, $\Sigma_{n\mathbf{k}} \in \mathbb{R}^{B \times B}$ are defined by

$$\pi_{\mathbf{k}} := \prod_{j=1}^{M} \pi_{jk_j}, \quad \mu_{n\mathbf{k}} := \sum_{j=1}^{M} \alpha_{nj} \mu_{jk_j},$$
 229

$$\mathbf{\Sigma}_{n\mathbf{k}} := \sum_{j=1}^{M} \alpha_{nj}^2 \mathbf{\Sigma}_{jk_j} + \mathbf{D}. \tag{6}$$

The proof is detailed using a characteristic function (c.f.) approach.

We first consider the distribution of the intermediate variable $\mathbf{z}_n = \sum_{j=1}^M \mathbf{m}_{nj} \alpha_{nj}$. The c.f. of $f_{\mathbf{m}_j}$ in (4), $\phi_{\mathbf{m}_j}$ (t) : $\mathbb{R}^B \to \mathbb{C}$, is given by

$$\phi_{\mathbf{m}_{j}}(\mathbf{t}) = \mathbb{E}_{\mathbf{m}_{j}}\left(e^{i\mathbf{t}^{T}\mathbf{x}}\right) = \int_{\mathbb{R}^{B}} e^{i\mathbf{t}^{T}\mathbf{x}} f_{\mathbf{m}_{j}}(\mathbf{x}) d\mathbf{x}$$

$$= \sum_{k=1}^{K_{j}} \pi_{jk} \int_{\mathbb{R}^{B}} e^{i\mathbf{t}^{T}\mathbf{x}} \mathcal{N}\left(\mathbf{x} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right) d\mathbf{x}$$

$$= \sum_{k=1}^{K_{j}} \pi_{jk} \phi_{jk}(\mathbf{t}), \tag{7}$$

where ϕ_{jk} (t) denotes the c.f. of the Gaussian distribution $\mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}_{jk},\boldsymbol{\Sigma}_{jk}\right)$ as

$$\phi_{jk}(\mathbf{t}) := \exp\left(i\mathbf{t}^T \boldsymbol{\mu}_{jk} - \frac{1}{2}\mathbf{t}^T \boldsymbol{\Sigma}_{jk}\mathbf{t}\right).$$
 (8)

Assuming $\mathbf{m}_{n1}, \dots, \mathbf{m}_{nM}$ are independent, we can obtain the c.f. of the linear combination of these \mathbf{m}_{nj} by multiplying (7) as

245
$$\phi_{\mathbf{z}_n}\left(\mathbf{t}
ight)$$

$$= \phi_{\mathbf{m}_{n1}\alpha_{n1} + \dots + \mathbf{m}_{nM}\alpha_{nM}}(\mathbf{t}) = \prod_{j=1}^{M} \phi_{\mathbf{m}_{j}}(\alpha_{nj}\mathbf{t})$$

$$= \sum_{k_1=1}^{K_1} \cdots \sum_{k_M=1}^{K_M} \pi_{1k_1} \cdots \pi_{Mk_M} \phi_{1k_1} (\alpha_{n1} \mathbf{t}) \cdots \phi_{Mk_M} (\alpha_{nM} \mathbf{t}).$$

Let K, \mathbf{k} , $\pi_{\mathbf{k}}$ be defined as in Theorem 1. We can write the above multiple summations in an elegant way:

$$\phi_{\mathbf{z}_n}(\mathbf{t}) = \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \phi_{n\mathbf{k}}(\mathbf{t}), \tag{9}$$

where $\pi_{\mathbf{k}} \geq 0$, $\sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} = 1$ and

$$\phi_{n\mathbf{k}}(\mathbf{t}) := \phi_{1k_1}(\alpha_{n1}\mathbf{t})\cdots\phi_{Mk_M}(\alpha_{nM}\mathbf{t})$$

$$= \exp \left\{ i \mathbf{t}^T \left(\sum_{j=1}^M \alpha_{nj} \boldsymbol{\mu}_{jk_j} \right) - \frac{1}{2} \mathbf{t}^T \left(\sum_{j=1}^M \alpha_{nj}^2 \boldsymbol{\Sigma}_{jk_j} \right) \mathbf{t} \right\},\,$$

where (8) is used. Since $\phi_{n\mathbf{k}}(\mathbf{t})$ also has a form of c.f. of a Gaussian distribution, the corresponding distribution turns out to be $\mathcal{N}\left(\mathbf{x}|\sum_{j}\alpha_{nj}\boldsymbol{\mu}_{jk_{j}},\sum_{j}\alpha_{nj}^{2}\boldsymbol{\Sigma}_{jk_{j}}\right)$. Hence, the distribution of \mathbf{z}_{n} can be obtained by the Fourier transform of (9)

$$f_{\mathbf{z}_{n}}(\mathbf{z}_{n}) = \frac{1}{(2\pi)^{B}} \int_{\mathbb{R}^{B}} e^{-i\mathbf{t}^{T}\mathbf{z}_{n}} \phi_{\mathbf{z}_{n}}(\mathbf{t}) d\mathbf{t}$$

$$= \frac{1}{(2\pi)^{B}} \int_{\mathbb{R}^{B}} e^{-i\mathbf{t}^{T}\mathbf{z}_{n}} \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \phi_{n\mathbf{k}}(\mathbf{t}) d\mathbf{t}$$

$$= \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \mathcal{N} \left(\mathbf{z}_{n} | \sum_{i=1}^{M} \alpha_{nj} \boldsymbol{\mu}_{jk_{j}}, \sum_{i=1}^{M} \alpha_{nj}^{2} \boldsymbol{\Sigma}_{jk_{j}} \right), \quad (10)$$

which is still a mixture of Gaussians.

After finding the distribution of the linear combination, we can add the noise term to find the distribution of y_n . Suppose the noise also follows a Gaussian distribution,

 $p(\mathbf{n}_n) := f_{\mathbf{n}_n}(\mathbf{n}_n) = \mathcal{N}(\mathbf{n}_n|\mathbf{0},\mathbf{D})$, where \mathbf{D} is the noise covariance matrix. We assume that the noise at different wavelengths is independent (σ_k^2) being the noise variance of the kth band), i.e. $\mathbf{D} = \mathrm{diag}\left(\sigma_1^2, \sigma_2^2, \ldots, \sigma_B^2\right) \in \mathbb{R}^{B \times B}$ (if it is not independent, the noise can actually be easily whitened to be independent as in [36]). Its c.f. has the following form

$$\phi_{\mathbf{n}_n}(\mathbf{t}) = \exp\left(-\frac{1}{2}\mathbf{t}^T\mathbf{D}\mathbf{t}\right) \tag{11}$$

by (8). Then the c.f. of \mathbf{y}_n can be obtained by multiplying (9) and (11) (as \mathbf{z}_n and \mathbf{n}_n are independent)

$$\phi_{\mathbf{y}_{n}}\left(\mathbf{t}\right)=\phi_{\mathbf{z}_{n}}\left(\mathbf{t}\right)\phi_{\mathbf{n}_{n}}\left(\mathbf{t}\right)=\sum_{\mathbf{k}\in\mathcal{K}}\pi_{\mathbf{k}}\phi_{\mathbf{n}_{n}}\left(\mathbf{t}\right)\phi_{n\mathbf{k}}\left(\mathbf{t}\right)$$
 274

$$= \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \exp \left\{ i \mathbf{t}^T \boldsymbol{\mu}_{n\mathbf{k}} - \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma}_{n\mathbf{k}} \mathbf{t} \right\},$$

where $\mu_{n\mathbf{k}}$ and $\Sigma_{n\mathbf{k}}$ are defined in (6). Finally, the distribution of \mathbf{y} can be shown to be (5) by the Fourier transform again as in (10).

If $K = \{1\} \times \{1\} \times \dots \times \{1\}$, i.e. each endmember has only one Gaussian component, we have $\pi_{11} = 1, \dots, \pi_{M1} = 1$, then $\pi_{\mathbf{k}} = \pi_{11} \cdots \pi_{M1} = 1$. The distribution of \mathbf{y}_n becomes

$$p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right) = \mathcal{N}\left(\mathbf{y}_{n}|\sum_{j=1}^{M}a_{nj}\boldsymbol{\mu}_{j1},\sum_{j=1}^{M}\alpha_{nj}^{2}\boldsymbol{\Sigma}_{j1} + \mathbf{D}\right),$$
(12) 283

which is exactly the NCM in (3).

B. Another Perspective

Theorem 1 obtains the density of each pixel by directly performing a r.v.t. based on the LMM, which can be used to estimate the abundances and distribution parameters. Here, we will obtain the density from another perspective, which provides a foundation to estimate the endmembers for each pixel. Again, let the noise follow the density function $p(\mathbf{n}_n) := \mathcal{N}(\mathbf{n}_n | \mathbf{0}, \mathbf{D})$. Considering $\{\mathbf{m}_{nj}\}$ and $\{\alpha_{nj}\}$ as fixed values, the r.v.t. $\mathbf{y}_n = \sum_j \mathbf{m}_{nj} \alpha_{nj} + \mathbf{n}_n$ implies that the density of \mathbf{y}_n is given by

$$p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\mathbf{M}_{n},\mathbf{D}\right) = \mathcal{N}\left(\mathbf{y}_{n}|\sum_{j}\mathbf{m}_{nj}\alpha_{nj},\mathbf{D}\right)$$
(13) 295

where $\mathbf{M}_n = [\mathbf{m}_{n1}, \dots, \mathbf{m}_{nM}]^T \in \mathbb{R}^{M \times B}$ are the endmembers for the *n*th pixel. We have the following theorem which gives the same result as in Theorem 1.

Theorem 2: If the random variables $\{\mathbf{m}_{nj}: j=1,\ldots,M\}$ follow GMM distributions

$$p\left(\mathbf{m}_{nj}|\mathbf{\Theta}\right) := \sum_{k=1}^{K_j} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj}|\boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right),$$
 301

and they are independent, i.e.

$$p\left(\mathbf{M}_{n}|\mathbf{\Theta}\right) = \prod_{i=1}^{M} p\left(\mathbf{m}_{nj}|\mathbf{\Theta}\right),\tag{14}$$

352

355

357

366

371

372

376

378

380

TABLE I
VALUES FOR THE VARIOUS QUANTITIES IN THE SIMPLE EXAMPLE

k	$\pi_{\mathbf{k}}$	$oldsymbol{\mu_{n\mathbf{k}}}$ in (6)
(1, 1, 1, 1)	0.06	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{21} + \alpha_{n3}\boldsymbol{\mu}_{31} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1, 2, 1, 1)	0.14	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{22} + \alpha_{n3}\boldsymbol{\mu}_{31} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1, 1, 2, 1)	0.12	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{21} + \alpha_{n3}\boldsymbol{\mu}_{32} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1, 2, 2, 1)	0.28	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{22} + \alpha_{n3}\boldsymbol{\mu}_{32} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1, 1, 3, 1)	0.12	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{21} + \alpha_{n3}\boldsymbol{\mu}_{33} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1, 2, 3, 1)	0.28	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{22} + \alpha_{n3}\boldsymbol{\mu}_{33} + \alpha_{n4}\boldsymbol{\mu}_{41}$

then the conditional density $p(\mathbf{y}_n|\boldsymbol{\alpha}_n,\boldsymbol{\Theta},\mathbf{D})$ obtained by marginalizing \mathbf{M}_n in $p(\mathbf{y}_n,\mathbf{M}_n|\boldsymbol{\alpha}_n,\boldsymbol{\Theta},\mathbf{D})$ has the same form as in Theorem 1:

$$p(\mathbf{y}_n|\boldsymbol{\alpha}_n, \boldsymbol{\Theta}, \mathbf{D}) = \int p(\mathbf{y}_n|\boldsymbol{\alpha}_n, \mathbf{M}_n, \mathbf{D}) p(\mathbf{M}_n|\boldsymbol{\Theta}) d\mathbf{M}_n$$
$$= \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \mathcal{N}(\mathbf{y}_n|\boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}}),$$

where
$$p(\mathbf{y}_n|\boldsymbol{\alpha}_n, \mathbf{M}_n, \mathbf{D}) = \mathcal{N}(\mathbf{y}_n|\sum_j \mathbf{m}_{nj}\alpha_{nj}, \mathbf{D}).$$

The proof is much more complicated (in terms of algebra) and therefore relegated to the supplemental material of the paper.

C. An Example

306

308

311

313

314

315

316

317

322

324

326

329

330

331

334

335

336

337

338

We give an example to illustrate the basic idea of this paper. Suppose we have M = 4 endmembers with $K_1 = 1$, $K_2 = 2$, $K_3 = 3$, $K_4 = 1$. Their distributions follow (4) with $\mu_{ik}, \Sigma_{ik}, j = 1, 2, 3, 4, k = 1, ..., K_i$. Let the weights of these components be $\pi_{11} = \pi_{41} = 1$, $\pi_{21} = 0.3$, $\pi_{22} = 0.7$, $\pi_{31} = 0.2$, $\pi_{32} = 0.4$, $\pi_{33} = 0.4$. Then, K has 6 entries from the Cartesian product, $\{1\} \times \{1, 2\} \times \{1, 2, 3\} \times \{1\}$. We list the values for $\pi_{\mathbf{k}}$, $\mu_{n\mathbf{k}}$ in Table I. For example, for $\mathbf{k} = (1, 2, 3, 1), \ \pi_{\mathbf{k}} = \pi_{11}\pi_{22}\pi_{33}\pi_{41} = 0.28.$ The value of $\mu_{n\mathbf{k}}$ is a linear combination of μ_{ik} (pick one component for each j) based on the configuration k. Hence, the distribution of y_n in (5) is a Gaussian mixture of 6 components with π_k , $\mu_{n\mathbf{k}}$ given in Table I ($\Sigma_{n\mathbf{k}}$ can be derived similar to $\mu_{n\mathbf{k}}$). Recalling the intuition in Fig. 2, we will show that applying it to hyperspectral unmixing will force each pixel to match all the $\mu_{n\mathbf{k}}$ s, but with emphasis determined by $\pi_{n\mathbf{k}}$.

III. GAUSSIAN MIXTURE MODEL FOR ENDMEMBER VARIABILITY

A. The GMM for Hyperspectral Unmixing

Based on the analysis in Section II, we can model the conditional distribution of all the pixels $\mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_N]^T \in \mathbb{R}^{N \times B}$ given all the abundances $\mathbf{A} := [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_N]^T \in \mathbb{R}^{N \times M}$ ($\boldsymbol{\alpha}_n := [\alpha_{n1}, \dots, \alpha_{nM}]^T$) and GMM parameters, which leads to a maximum *a posteriori* (MAP) problem. Using the result in (5) and assuming the conditional distributions of \mathbf{y}_n are independent, the distribution of \mathbf{Y} given $\mathbf{A}, \boldsymbol{\Theta}, \mathbf{D}$ becomes

$$p(\mathbf{Y}|\mathbf{A}, \mathbf{\Theta}, \mathbf{D}) = \prod_{n=1}^{N} p(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n}, \mathbf{\Theta}, \mathbf{D}).$$
 (15)

Based on the hyperspectral unmixing context, we can set the priors for A. Suppose we use the same prior on A as in [37], i.e.

$$p(\mathbf{A}) \propto \exp\left\{-\frac{\beta_1}{2} \operatorname{Tr}\left(\mathbf{A}^T \mathbf{L} \mathbf{A}\right) + \frac{\beta_2}{2} \operatorname{Tr}\left(\mathbf{A}^T \mathbf{A}\right)\right\}$$
$$= \exp\left\{-\frac{\beta_1}{2} \operatorname{Tr}\left(\mathbf{A}^T \mathbf{K} \mathbf{A}\right)\right\}, \tag{16}$$

where **L** is a *graph Laplacian* matrix constructed from w_{nm} , $n, m = 1, \ldots, N$ with $w_{nm} = e^{-\|\mathbf{y}_n - \mathbf{y}_m\|^2/2B\eta^2}$ for neighboring pixels and 0 otherwise. We have $\operatorname{Tr}(\mathbf{A}^T\mathbf{L}\mathbf{A}) = \frac{1}{2}\sum_{n,m}w_{nm}\|\boldsymbol{\alpha}_n - \boldsymbol{\alpha}_m\|^2 = \mathbf{L} - \frac{\beta_2}{\beta_1}\mathbf{I}_N$ (suppose $\beta_1 \neq 0$) with β_1 controlling smoothness and β_2 controlling sparsity of the abundance maps.

From the conditional density function and the priors, Bayes' theorem says the posterior is given by

$$p(\mathbf{A}, \mathbf{\Theta}|\mathbf{Y}, \mathbf{D}) \propto p(\mathbf{Y}|\mathbf{A}, \mathbf{\Theta}, \mathbf{D}) p(\mathbf{A}) p(\mathbf{\Theta}),$$
 (17)

where $p(\Theta)$ is assumed to follow a uniform distribution. Maximizing $p(\mathbf{A}, \Theta | \mathbf{Y}, \mathbf{D})$ is equivalent to minimizing $-\log p(\mathbf{A}, \Theta | \mathbf{Y}, \mathbf{D})$, which reduces to the following form by combining (5), (15), (16) and (17):

$$\mathcal{E}\left(\mathbf{A},\mathbf{\Theta}\right) = -\sum_{n=1}^{N}\log\sum_{\mathbf{k}\in\mathcal{K}}\pi_{\mathbf{k}}\mathcal{N}\left(\mathbf{y}_{n}|\boldsymbol{\mu}_{n\mathbf{k}},\boldsymbol{\Sigma}_{n\mathbf{k}}\right) + \mathcal{E}_{\mathrm{prior}}(\mathbf{A}),$$

s.t.
$$\pi_{\mathbf{k}} \ge 0$$
, $\sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} = 1$, $\alpha_{nj} \ge 0$, $\sum_{j=1}^{M} \alpha_{nj} = 1$, $\forall n$ (18)

where $\mathcal{E}_{prior}(\mathbf{A}) = \frac{\beta_1}{2} \text{Tr} (\mathbf{A}^T \mathbf{K} \mathbf{A})$, and $\boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}}$ are defined in (6).

B. Relationships to Least-Squares, NCM, and MESMA

Let us focus on the first term in (18) and call it the *data fidelity term*. We can relate it to NCM and the least-squares term $\sum_n \|\mathbf{y}_n - \sum_j \alpha_{nj} \mathbf{m}_j\|^2$ as used in previous research. The data fidelity term in NCM follows (3) and is based on minimizing the negative log-likelihood

$$-\log p\left(\mathbf{Y}\right) = -\log \prod_{n=1}^{N} p\left(\mathbf{y}_{n}\right) = -\sum_{n=1}^{N} \log \mathcal{N}\left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n1}, \boldsymbol{\Sigma}_{n1}\right)$$

$$\tag{19}$$

by assuming \mathbf{y}_n s are independent, where $\boldsymbol{\mu}_{n1} := \sum_j \alpha_{nj} \boldsymbol{\mu}_j$, $\boldsymbol{\Sigma}_{n1} := \sum_j \alpha_{nj}^2 \boldsymbol{\Sigma}_j + \sigma^2 \mathbf{I}_B$. Expanding (19) using the form of the Gaussian distribution leads to the objective function

$$\sum_{n=1}^{N} \log |\mathbf{\Sigma}_{n1}| + \sum_{n=1}^{N} (\mathbf{y}_{n} - \boldsymbol{\mu}_{n1})^{T} \, \boldsymbol{\Sigma}_{n1}^{-1} (\mathbf{y}_{n} - \boldsymbol{\mu}_{n1}). \quad (20)$$

We can see that the least-squares minimization is a special case of NCM with $\|\mathbf{\Sigma}_j\|_F \to 0$, i.e. when there is little endmember variability.

The proposed GMM further generalizes NCM from a statistical perspective. Since π_{jk} represents the prior probability of the latent variable in a GMM, $\pi_{\mathbf{k}}$ represents the prior probability of picking a combination. If we see \mathbf{k} as a (discrete) random variable whose sample space is \mathcal{K} , (5) can be seen as

$$p(\mathbf{y}_n|\mathbf{\alpha}_n, \mathbf{\Theta}, \mathbf{D}) = \sum_{\mathbf{k} \in \mathcal{K}} p(\mathbf{k}) p(\mathbf{y}_n|\mathbf{k}, \mathbf{\alpha}_n, \mathbf{\Theta}, \mathbf{D}),$$
 38

435

436

446

454

455

458

460

462

464

465

466

468

469

383

385

387

389

390

391

393

394

395

396

398

399

400

402

403

404

406

408

409

410

411

412

413

414

415

416

417

418

420

421

422

424

425

427

428

431

433

where $p(\mathbf{k}) = \pi_{\mathbf{k}}$ and $p(\mathbf{y}_n | \mathbf{k}, \boldsymbol{\alpha}_n, \boldsymbol{\Theta}, \mathbf{D}) = \mathcal{N}(\mathbf{y}_n | \mathbf{k})$ $\mu_{n\mathbf{k}}, \Sigma_{n\mathbf{k}}$). From this perspective, each pixel is generated by first sampling k, then sampling a Gaussian distribution determined by \mathbf{k} , $\mathbf{\Theta}$. Unlike NCM that tries to make each \mathbf{y}_n close to μ_{n1} which is a linear combination of a fixed set $\{\mu_i\}$, GMM further generalizes it by trying to make y_n close to every $\mu_{n\mathbf{k}}$ which are all the possible linear combinations of $\{\mu_{ik}\}$. It makes sense that the summation in (18) is weighted by $\pi_{\mathbf{k}}$ in a way that if one combination has a high probability to appear, i.e. $\pi_{\mathbf{k}}$ is larger for a certain \mathbf{k} , the effort is biased to make \mathbf{y}_n closer to this particular $\boldsymbol{\mu}_{n\mathbf{k}}$. Fig. 2 shows the differences among these.

The widely adopted MESMA takes a library of endmember spectra as input, tries all the combinations and pick the combination with least reconstruction error. The philosophy is similar to our model despite the fundamental difference that MESMA is explicit whereas we are implicit in terms of linear combinations. Compared to MESMA, the GMM approach separates the library into M groups where each group represents a material and is clustered into several centers, such that the combination can only take place by picking one center from each group. Also, the size of each cluster affects the probability of picking its center. Hence, our model can adapt to very large library sizes as long as the number of clusters does not increase too much.

C. Optimization

Estimating the parameters of GMMs has been studied extensively, from early expectation maximization (EM) from the statistical community to projection based clustering from the computer science community [38], [39]. There are simple and deterministic algorithms, which usually require the centers of Gaussian be separable. However, we face a more challenging problem since each pixel is generated by a different GMM determined by the coefficients α_n . Since EM can be seen as a special case of Majoriziation-Minimization algorithms [40], which is more flexible, we adopt this approach. Considering that we have too many parameters A, Θ to update in the M step, they are updated sequentially as long as the complete data log-likelihood increases. This is also called generalized expectation maximization (GEM) [41].

Following the routine of EM, the E step calculates the posterior probability of the latent variable given the observed data and old parameters

$$\gamma_{n\mathbf{k}} = \frac{\pi_{\mathbf{k}} \mathcal{N} \left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}} \right)}{\sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \mathcal{N} \left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}} \right)}.$$
 (21)

The M step usually maximizes the expected value of the complete data log-likelihood. Here, we have priors in the Bayesian formulation. Hence, we need to minimize

$$\mathcal{E}_{M} = -\sum_{n=1}^{N} \sum_{\mathbf{k} \in \mathcal{K}} \gamma_{n\mathbf{k}} \left\{ \log \pi_{\mathbf{k}} + \log \mathcal{N} \left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}} \right) \right\} + \mathcal{E}_{\text{prior}}.$$

This leads to a common update step for π_k as

$$\pi_{\mathbf{k}} = \frac{1}{N} \sum_{n=1}^{N} \gamma_{n\mathbf{k}}.$$
 (23)

We now focus on updating $\{\mu_{jk}, \Sigma_{jk}\}$ and **A**. To achieve this, we require the derivatives of \mathcal{E}_M in (22) w.r.t. $\boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{jk}, \alpha_{nj}$. After some tedious algebra using (6), we get

$$\frac{\partial \mathcal{E}_M}{\partial \boldsymbol{\mu}_{jl}} = -\sum_{n=1}^N \sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \alpha_{nj} \boldsymbol{\lambda}_{n\mathbf{k}}$$
 (24)

$$\frac{\partial \mathcal{E}_M}{\partial \mathbf{\Sigma}_{jl}} = -\sum_{n=1}^N \sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \alpha_{nj}^2 \mathbf{\Psi}_{n\mathbf{k}},\tag{25}$$

$$\frac{\partial \mathcal{E}_{M}}{\partial \alpha_{nj}} = -\sum_{\mathbf{k} \in \mathcal{K}} \boldsymbol{\lambda}_{n\mathbf{k}}^{T} \boldsymbol{\mu}_{jk_{j}} - 2\alpha_{nj} \sum_{\mathbf{k} \in \mathcal{K}} \operatorname{Tr} \left(\boldsymbol{\Psi}_{n\mathbf{k}}^{T} \boldsymbol{\Sigma}_{jk_{j}} \right) + \beta_{1} \left(\mathbf{K} \mathbf{A} \right)_{nj}, \tag{26}$$

where $\lambda_{n\mathbf{k}} \in \mathbb{R}^{B \times 1}$ and $\Psi_{n\mathbf{k}} \in \mathbb{R}^{B \times B}$ are given by

$$\lambda_{n\mathbf{k}} = \gamma_{n\mathbf{k}} \mathbf{\Sigma}_{n\mathbf{k}}^{-1} (\mathbf{y}_n - \boldsymbol{\mu}_{n\mathbf{k}}), \tag{27}$$

$$\Psi_{n\mathbf{k}} = \frac{1}{2} \gamma_{n\mathbf{k}} \mathbf{\Sigma}_{n\mathbf{k}}^{-T} \left(\mathbf{y}_n - \boldsymbol{\mu}_{n\mathbf{k}} \right) \left(\mathbf{y}_n - \boldsymbol{\mu}_{n\mathbf{k}} \right)^T \mathbf{\Sigma}_{n\mathbf{k}}^{-T} - \frac{1}{2} \gamma_{n\mathbf{k}} \mathbf{\Sigma}_{n\mathbf{k}}^{-T}.$$
(28)

It is better to represent the derivatives in matrix forms for the sake of implementation convenience. Considering the multiple summations in (24), (25) and (26), we can write them as

$$\frac{\partial \mathcal{E}_M}{\partial \mu_{jl}} = -\sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \left(\mathbf{A}^T \mathbf{\Lambda}_{\mathbf{k}} \right)_j, \tag{29}$$

$$\frac{\partial \mathcal{E}_M}{\partial \text{vec}\left(\mathbf{\Sigma}_{jl}\right)} = -\sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \left((\mathbf{A} \circ \mathbf{A})^T \, \mathbf{\Psi}_{\mathbf{k}} \right)_j, \tag{30}$$

$$\frac{\partial \mathcal{E}_M}{\partial \mathbf{A}} = -\sum_{\mathbf{k} \in \mathcal{K}} \mathbf{\Lambda}_{\mathbf{k}} \mathbf{R}_{\mathbf{k}}^T - 2\mathbf{A} \circ \sum_{\mathbf{k} \in \mathcal{K}} \mathbf{\Psi}_{\mathbf{k}} \mathbf{S}_{\mathbf{k}}^T + \beta_1 \mathbf{K} \mathbf{A}, \quad (31)$$

where $\mathbf{\Lambda}_{\mathbf{k}} \in \mathbb{R}^{N \times B}$, $\mathbf{\Psi}_{\mathbf{k}} \in \mathbb{R}^{N \times B^2}$ denote the matrices formed by $\{\lambda_{n\mathbf{k}}, \Psi_{n\mathbf{k}}\}$ as follows

$$\Lambda_{\mathbf{k}} := [\lambda_{1\mathbf{k}}, \lambda_{2\mathbf{k}}, \dots, \lambda_{N\mathbf{k}}]^{T},$$

$$\Psi_{\mathbf{k}} := [\operatorname{vec}(\Psi_{1\mathbf{k}}), \operatorname{vec}(\Psi_{2\mathbf{k}}), \dots, \operatorname{vec}(\Psi_{N\mathbf{k}})]^{T},$$
453

and
$$\mathbf{R}_{\mathbf{k}} \in \mathbb{R}^{M \times B}$$
, $\mathbf{S}_{\mathbf{k}} \in \mathbb{R}^{M \times B^2}$ are defined by

$$\mathbf{R}_{\mathbf{k}} := \begin{bmatrix} \mu_{1k_1}, \mu_{2k_2}, \dots, \mu_{Mk_M} \end{bmatrix}^T, \tag{32}$$

$$\mathbf{S}_{\mathbf{k}} := \left[\operatorname{vec} \left(\mathbf{\Sigma}_{1k_1} \right), \operatorname{vec} \left(\mathbf{\Sigma}_{2k_2} \right), \dots, \operatorname{vec} \left(\mathbf{\Sigma}_{Mk_M} \right) \right]^T.$$
 (33)

The minimum of \mathcal{E}_M corresponds to $\frac{\partial \mathcal{E}_M}{\partial \mu_{ij}} = 0$, $\frac{\partial \mathcal{E}_M}{\partial \Sigma_{ij}} = 0$, and $\frac{\partial \mathcal{E}_M}{\partial \mathbf{A}} = 0$ if the optimization problem is unconstrained. However, since we have the non-negativity and sum-to-one constraint to α_{nj} and positive definite constraint of Σ_{jk} , minimizing \mathcal{E}_M is very difficult. Therefore, in each M step, we only decrease this objective function by *projected gradient* descent (please see [42 and 43, Sec. 2.3]) using (29), (30) and (31), where the projection functions for **A** and $\{\Sigma_{ik}\}$ are the same as in [37].

Finally, from the estimated π_k , we can recover the sets of weights as $\pi_{il} = \sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_i} \pi_{\mathbf{k}}$.

D. Model Selection

(22)

The number of components K_i can be specified or estimated from the data. For the latter case, we have some pure pixels and estimate K_i by deploying a standard

model selection method. Suppose we have N_j pure pixels $\mathbf{Y}_j := \begin{bmatrix} \mathbf{y}_1^j, \mathbf{y}_2^j, \dots, \mathbf{y}_{N_j}^j \end{bmatrix}^T \in \mathbb{R}^{N_j \times B}$ for the jth endmember, $f_{\mathbf{m}_j}\left(\mathbf{y}|\mathbf{\Theta}_j\right)$ is the estimated density function with $\mathbf{\Theta}_j := \{\pi_{jk}, \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk} : k = 1, \dots, K_j\}$, $g_{\mathbf{m}_j}\left(\mathbf{y}\right)$ is the true density function. The information criterion based model selection approach tries to find K_j that minimizes their difference, e.g. the Kullback-Leibler (KL) divergence

$$\mathcal{D}_{\mathrm{KL}}\left(g_{\mathbf{m}_{j}} \| f_{\mathbf{m}_{j}}\right) = \int_{\mathbb{R}^{B}} g_{\mathbf{m}_{j}}\left(\mathbf{y}\right) \log \frac{g_{\mathbf{m}_{j}}\left(\mathbf{y}\right)}{f_{\mathbf{m}_{j}}\left(\mathbf{y} | \mathbf{\Theta}_{j}\right)} d\mathbf{y}$$

$$\approx -\frac{1}{N_{j}} \sum_{n=1}^{N_{j}} \log f_{\mathbf{m}_{j}}\left(\mathbf{y}_{n}^{j} | \mathbf{\Theta}_{j}\right) + \text{const},$$

where the approximation of $\int g_{\mathbf{m}_j}(\mathbf{y}) \log f_{\mathbf{m}_j}(\mathbf{y}|\mathbf{\Theta}_j) d\mathbf{y}$ by the log-likelihood is usually biased as the empirical distribution function is closer to the fitted distribution than the true one. Akaike's information criterion is one way to approximate the bias. Here, we use the cross-validation-based information criterion (CVIC) to correct for the bias [44], [45]. Let

$$\mathcal{L}_{\mathbf{Y}_{j}}\left(\mathbf{\Theta}_{j}\right) = \sum_{n=1}^{N_{j}} \log f_{\mathbf{m}_{j}}\left(\mathbf{y}_{n}^{j} | \mathbf{\Theta}_{j}\right). \tag{34}$$

The *V*-fold cross validation (we use V = 5 here) divides the input set \mathbf{Y}_j into V subsets $\left\{\mathbf{Y}_j^1, \mathbf{Y}_j^2, \dots, \mathbf{Y}_j^V\right\}$ with equal sizes. Then for each subset \mathbf{Y}_j^v , $v = 1, \dots, V$, the remaining data are used to replace \mathbf{Y}_j in (34) such that (34) is maximized by $\mathbf{\Theta}_j^v$. Then $\mathcal{L}_{K_j} = \sum_v \mathcal{L}_{\mathbf{Y}_j^v} \left(\mathbf{\Theta}_j^v\right)$ is evaluated and the optimal $\hat{K}_j = \arg\max_{K_j} \mathcal{L}_{K_j}$.

E. Implementation Details

The algorithm can be implemented in a supervised or unsupervised manner. In both cases, because of the large computational cost, we project the pixel data to a low dimensional space by principal component analysis (PCA) and perform the optimization, the result then projected back to the original space. Let $\mathbf{E} \in \mathbb{R}^{B \times d}$ be the projection matrix and $\mathbf{c} \in \mathbb{R}^{B}$ be the translation vector, then

$$\mathbf{E}^{T} (\mathbf{y}_{n} - \mathbf{c}) = \sum_{j=1}^{M} \mathbf{E}^{T} (\mathbf{m}_{nj} - \mathbf{c}) \alpha_{nj} + \mathbf{E}^{T} \mathbf{n}_{n}.$$

This means that for the projected pixels, the *j*th endmember $\mathbf{m}'_{nj} = \mathbf{E}^T (\mathbf{m}_{nj} - \mathbf{c})$ follows a distribution

$$p\left(\mathbf{m}_{nj}'|\mathbf{\Theta}\right) = \sum_{k=1}^{K_j} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj}'|\mathbf{E}^T \left(\boldsymbol{\mu}_{jk} - \mathbf{c}\right), \mathbf{E}^T \boldsymbol{\Sigma}_{jk} \mathbf{E}\right)$$

and the noise $\mathbf{n}'_n = \mathbf{E}^T \mathbf{n}_n$ follows $\mathcal{N}(\mathbf{n}'_n | \mathbf{0}, \mathbf{E}^T \mathbf{D} \mathbf{E})$.

In the supervised unmixing scenario, we assume that a library of endmember spectra is known. After estimating the number of components following Section III-D, and calculating Θ using the standard EM algorithm, we only need to update $\gamma_{n\mathbf{k}}$ by (21) and \mathbf{A} by (31) with $\pi_{\mathbf{k}}$, μ_{jk} and Σ_{jk} fixed. The initialization of \mathbf{A} can utilize the multiple combinations of means. For each α_n , we first set

 $\alpha_{n\mathbf{k}} \leftarrow (\mathbf{R}_{\mathbf{k}}\mathbf{R}_{\mathbf{k}}^T + \epsilon \mathbf{I}_M)^{-1} \mathbf{R}_{\mathbf{k}}\mathbf{y}_n$, then project it to the simplex space, and finally set $\alpha_n \leftarrow \alpha_n$ that minimizes the reconstruction error

In the unsupervised unmixing scenario, we will assume the resolution is high enough such that the hyperspectral image can be segmented into several regions where the interior pixels in each region are pure pixels. The optimization is performed in several steps, where we first obtain a segmentation result, then use CVIC to determine the number of components, and finally estimate $\bf A$ with $\bf \Theta$ fixed. The details are given as follows.

Step 1: Initialization. We start with $K_j = 1$, $\forall j$ and use K-means to find the initial means $\mathbf{R_1}$. The initial \mathbf{A} is set to $\mathbf{A} \leftarrow \mathbf{Y}\mathbf{R}_1^T \left(\mathbf{R}_1\mathbf{R}_1^T + \epsilon \mathbf{I}_M\right)^{-1}$ (by minimizing $\|\mathbf{Y} - \mathbf{A}\mathbf{R}_1\|_F^2$), then projected to the valid simplex space as in [37]. The initial covariance matrices are set to $\mathbf{\Sigma}_{j1} \leftarrow 0.1^2\mathbf{I}_B$, $\forall j$. For the noise matrix \mathbf{D} , although there is research focused on noise estimation [46], [47], endmember variability was not considered and validation was performed only for the simple LMM assumption. Hence, we use an empirical value $\mathbf{D} = 0.001^2\mathbf{I}_B$, which is usually much less than the variability of covariance matrices in (6).

Step 2: Segmentation. Given the initial conditions, we use the GEM algorithm to iteratively update $\gamma_{n\mathbf{k}}$ by (21), $\pi_{\mathbf{k}}$ by (23), μ_{jk} by (29), \mathbf{A} by (31) while keeping Σ_{jk} fixed. For $\gamma_{n\mathbf{k}}$ and $\pi_{\mathbf{k}}$, a direct update equation is available. For μ_{jk} , we can use gradient descent. For \mathbf{A} , since we have the non-negativity and sum-to-one constraints, a projected gradient descent similar to the one used in [37] can be applied. To ensure a segmentation effect, a large β_2 is used in this step.

Step 3: Model selection and abundance estimation. Using the segmentation-like abundance maps from the previous step, we can obtain the interior pixels \mathbf{Y}_j (assumed pure) by thresholding the abundances (e.g. $\alpha_{nj} > 0.99$) and performing image erosion to trim the boundaries with structure element size r_{se} (can be decreased gradually if large enough to trim all the pixels). Following Section III-D, we can determine the number of components K_j and further calculate $\mathbf{\Theta}_j$ by standard EM. Since β_2 is relatively large in the previous step, it is reduced by $\beta_2 \leftarrow \zeta \beta_2$ where $\zeta = 0.05$. Then we restart the optimization to estimate the abundances with $\mathbf{\Theta}$ fixed.

F. Complexity Analysis

The abundance estimation algorithm is an iterative process. Since we used projected gradient descent with adaptive step sizes, the number of iterations is usually not large as shown in [43] and [48]. For each iteration, it starts with calculating $\mu_{n\mathbf{k}}$ and $\Sigma_{n\mathbf{k}}$ in (6), where storing all $\mu_{n\mathbf{k}}$ ($\Sigma_{n\mathbf{k}}$) requires $O(|\mathcal{K}|NB)$ ($O(|\mathcal{K}|NB^2)$), the computation takes $O(|\mathcal{K}|NMB)$ ($O(|\mathcal{K}|NMB^2)$). Suppose the Cholesky factorization and the matrix inversion of a B by B matrix both take $O(B^3)$ time, and $N\gg B>M$. Evaluating $\log \mathcal{N}(\mathbf{y}_n|\mu_{n\mathbf{k}},\Sigma_{n\mathbf{k}})$ by the Cholesky factorization will take $O(B^3)$, hence updating all the $\gamma_{n\mathbf{k}}$ takes $O(|\mathcal{K}|NB^3)$, which is also the required time for evaluating the objective function (18). The calculation of $\lambda_{n\mathbf{k}}$, $\Psi_{n\mathbf{k}}$ (in (27) and (28)) will be

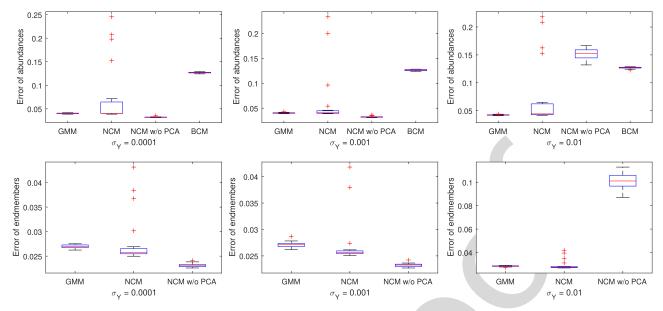


Fig. 3. Abundance and endmember error statistics from 20 synthetic images for each noise level in the supervised unmixing scenario.

dominated by the inversion of $\Sigma_{n\mathbf{k}}$ which takes $O\left(B^3\right)$, hence the overall calculation takes $O\left(|\mathcal{K}|\,NB^3\right)$ with storage the same as $\mu_{n\mathbf{k}}$ and $\Sigma_{n\mathbf{k}}$. Then if we move to calculating the derivatives in (29), (30) and (31), it is easy to verify that the computational costs are $O\left(|\mathcal{K}|\,NMB\right)$, $O\left(|\mathcal{K}|\,NMB^2\right)$, $O\left(|\mathcal{K}|\,NMB^2\right)$ respectively (Note that \mathbf{K} is a banded matrix so the computation involving it is linear). Reviewing the above process, we conclude that the spatial complexity is dominated by $O\left(|\mathcal{K}|\,NB^2\right)$ and the time complexity is dominated by $O\left(|\mathcal{K}|\,NB^2\right)$.

G. Estimation of Endmembers for Each Pixel

While the previous sections discuss the estimation of the abundances and endmember distribution parameters, they do not actually estimate the endmembers $\{\mathbf{m}_{nj}: n=1,\ldots,N, j=1,\ldots,M\}$ for each pixel. In this Section, we will discuss this additional problem and note its absence in the previous NCM literature.

Theorem 2 implies that we can view the proposed conditional density (5) as modeling the noise as a Gaussian random variable followed by marginalizing over \mathbf{M}_n , which is usually achieved by the evidence approximation in the machine learning literature due to the intractability of the integral ([49, Sec. 3.5]). Since we have $\mathbf{A}, \boldsymbol{\Theta}$ obtained from the previous Sections, we can get the posterior of \mathbf{M}_n from this model:

$$p\left(\mathbf{M}_{n}|\mathbf{y}_{n},\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right) \propto p\left(\mathbf{y}_{n},\mathbf{M}_{n}|\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right)$$
$$= p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\mathbf{M}_{n},\mathbf{D}\right) p\left(\mathbf{M}_{n}|\boldsymbol{\Theta}\right). \quad (35)$$

Maximizing $\log p(\mathbf{M}_n|\mathbf{y}_n, \boldsymbol{\alpha}_n, \boldsymbol{\Theta}, \mathbf{D})$ gives us another minimization problem

$$\mathcal{E}\left(\mathbf{M}_{n}\right) = \frac{1}{2} \left(\mathbf{y}_{n} - \mathbf{M}_{n}^{T} \boldsymbol{\alpha}_{n}\right)^{T} \mathbf{D}^{-1} \left(\mathbf{y}_{n} - \mathbf{M}_{n}^{T} \boldsymbol{\alpha}_{n}\right)$$
$$- \sum_{i=1}^{M} \log \sum_{k=1}^{K_{j}} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right) \quad (36)$$

obtained by plugging (13) and (14) into (35). Note that this objective function has an intuitive interpretation as the first term minimizes the reconstruction error while the second term forces the endmembers close to the centers of each GMM. The weight factor between the two terms is the noise. From an algebraic perspective, since there are also logarithms of sums of Gaussian functions in this objective, we can also use the EM algorithm for ease of optimization. In the E step, the soft membership is calculated by

$$\gamma_{njk} = \frac{\pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right)}{\sum_{k} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right)}, \quad k = 1, \dots, K_{j}.$$

In the M step, the derivative w.r.t. \mathbf{m}_{nj} is obtained as

$$\frac{\partial \mathcal{E}}{\partial \mathbf{m}_{nj}} = -\mathbf{D}^{-1} \left(\mathbf{y}_n - \mathbf{M}_n^T \boldsymbol{\alpha}_n \right) \alpha_{nj}$$

$$+ \sum_{k=1}^{K_j} \gamma_{njk} \boldsymbol{\Sigma}_{jk}^{-1} \left(\mathbf{m}_{nj} - \boldsymbol{\mu}_{jk} \right).$$
613

t descent in the M step for

Instead of deploying gradient descent in the M step for estimating the abundances, combining the derivatives for all j actually leads to a closed form solution

$$\operatorname{vec}\left(\mathbf{M}_{n}^{T}\right) = \left\{\boldsymbol{\alpha}_{n}\boldsymbol{\alpha}_{n}^{T} \otimes \mathbf{D}^{-1} + \operatorname{diag}\left(\mathbf{C}_{n1}, \dots, \mathbf{C}_{nM}\right)\right\}^{-1}$$

$$\left\{\operatorname{vec}\left(\mathbf{D}^{-1}\mathbf{y}_{n}\boldsymbol{\alpha}_{n}^{T}\right) + \mathbf{d}_{n}\right\}$$
618

where $\mathbf{C}_{nj} \in \mathbb{R}^{B \times B}$ and $\mathbf{d}_n := (\mathbf{d}_{n1}^T, \dots, \mathbf{d}_{nM}^T)^T \in \mathbb{R}^{MB \times 1}$ are defined as

$$\mathbf{C}_{nj} := \sum_{k=1}^{K_j} \gamma_{njk} \mathbf{\Sigma}_{jk}^{-1}, \ \mathbf{d}_{nj} := \sum_{k=1}^{K_j} \gamma_{njk} \mathbf{\Sigma}_{jk}^{-1} \boldsymbol{\mu}_{jk}.$$

In practice, despite the need to estimate a large $M \times B \times N$ tensor, the time cost is actually much less than the estimation of abundances because of the closed form update equation in the M step. An interesting fact is that γ_{njk} measures the closeness of estimated endmembers to clusters centers, hence

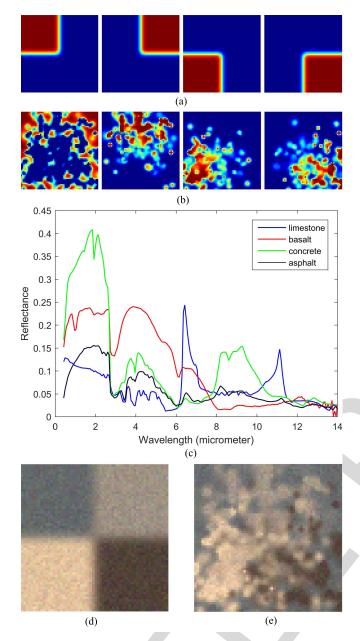


Fig. 4. Unsupervised synthetic dataset. (a) and (b) are abundance maps for two images. (c) shows original spectra from the ASTER library. (d) and (e) show the color images.

may provide a clue on which cluster is sampled to generate an endmember.

IV. RESULTS

In the following experiments, we implemented the algorithm in MATLAB® and compared the proposed GMM with NCM, BCM (spectral version with quadratic programming) [15] on synthetic and real images. As mentioned previously, for GMM, the original image data were projected to a subspace with 10 dimensions to speed up the computation for abundance estimation. NCM was implemented as a supervised algorithm wherein we input the ground truth pure pixels (in the image

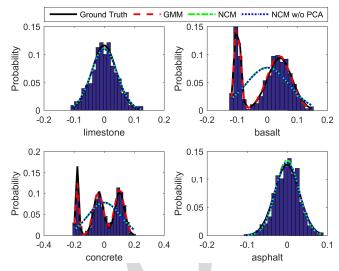


Fig. 5. Histograms of pure pixels for the 4 materials (when projected to a 1-dimensional space determined by performing PCA on the pure pixels of each material) and the ground truth and estimated distributions (also projected to the same direction) for the first image of the unsupervised synthetic dataset. The probability of each distribution is calculated by multiplying the value of the density function at each bin location with the bin size.

TABLE II $L_2 \ \, \text{DISTANCE Between the Fitted Distributions (GMM, NCM)} \\ \text{AND THE GROUND TRUTH DISTRIBUTIONS FOR THE FIRST} \\ \text{IMAGE OF THE UNSUPERVISED SYNTHETIC DATASET}$

$\times 10^{6}$	Limestone	Basalt	Concrete	Asphalt	Mean
GMM	4.45	3.46	3.41	4.28	3.85
NCM	4.27	5.86	4.95	4.02	4.77

TABLE III
ABUNDANCE ERRORS FOR THE UNSUPERVISED SYNTHETIC DATASET

	$\times 10^{-4}$	GMM	NCM	NCM w/o PCA	BCM
	Limestone	50	107	92	126
e 1	Basalt	40	74	67	158
Image	Concrete	41	66	62	186
ᄪ	Asphalt	69	141	123	292
	Mean	59	97	86	190
	Limestone	157	1086	396	231
e 2	Basalt	126	445	270	204
Image	Concrete	103	985	229	206
<u> </u>	Asphalt	225	170	706	445
	Mean	153	671	400	272

with extreme abundances), modeled them by Gaussian distributions, and obtained the abundance maps by maximizing the log-likelihood. We considered two versions of NCM, one in the same subspace as GMM (referred to as NCM), the other in the original spectral space (referred to as NCM without PCA). Since BCM is also a supervised unmixing algorithm, ground truth pure pixels were again taken as input and the results were the abundance maps. For GMM and the two versions of NCM, using the algorithm in Section III-G we can obtain the endmembers for each pixel. All the parameters of GMM (except the structure element size r_{se}) were set to $\beta_1 = 5$, $\beta_2 = 5$ unless specified throughout the experiments.

For comparison of endmember distributions, we calculated the L_2 distance $\left(\int |f(\mathbf{x}) - g(\mathbf{x})|^2 d\mathbf{x}\right)^{1/2}$ between the fitted distribution and the ground truth one, where the latter was only

¹The code of GMM is available on GitHub (https://github.com/zhouyuanzxcv/Hyperspectral).

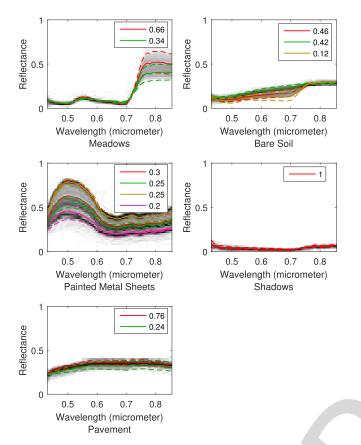


Fig. 6. Estimated GMM in the wavelength-reflectance space for the Pavia University dataset. The background gray image represents the histogram created by placing the pure pixel spectra into the reflectance bins at each wavelength. The different colors represent different components, where the solid curve is the center μ_{jk} , the dashed curves are $\mu_{jk} \pm 2\sigma_{jk} \mathbf{v}_{jk}$ (σ_{jk} is the square root of the large eigenvalue of Σ_{jk} while \mathbf{v}_{jk} is the corresponding eigenvector), and the legend shows the prior probabilities.

available for the synthetic dataset. For comparison of abundances, we calculated the root mean squared error (RMSE) $\left(\frac{1}{N}\sum_{n}|\alpha_{nj}^{GT}-\alpha_{nj}^{est}|^2\right)^{1/2}$ where α_{nj}^{GT} are the ground truth abundances and α_{nj}^{est} are the estimated values. Since only some pure pixels were identified as ground truth in the real datasets, we calculated error $j=\left(\frac{1}{|\mathcal{I}|}\sum_{n\in\mathcal{I}}|\alpha_{nj}^{GT}-\alpha_{nj}^{est}|^2\right)^{1/2}$ given the pure pixel index set \mathcal{I} . For comparison of endmembers, the same error formula and overall schema were used, i.e. for an index set \mathcal{I}_j of pure pixels for the jth endmember (in the real datasets), error $j=\frac{1}{|\mathcal{I}_j|}\sum_{n\in\mathcal{I}_j}\left(\frac{1}{B}\|\mathbf{m}_{nj}^{GT}-\mathbf{m}_{nj}^{est}\|^2\right)^{1/2}$.

A. Synthetic Datasets

The algorithms were tested for two cases of synthetic images, a supervised case and an unsupervised case.

1) Supervised: In this case, a library of ground truth endmembers were input and the abundances were estimated. The images were of size 60×60 with 103 wavelengths from 430 nm to 860 nm (\leq 5 nm spectral resolution) and created with two endmember classes, meadows and painted metal sheets, whose spectra were drawn randomly from the ground truth of the Pavia University dataset (shown in Fig. 1, meadows have 309 samples and painted metal sheets have

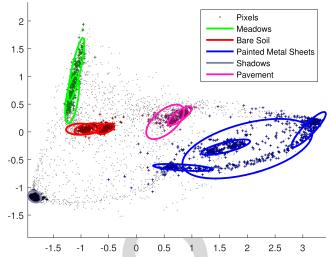


Fig. 7. Scatter plot of the Pavia University dataset with the estimated GMM. The gray dots are the projected pixels by PCA. The darkened dots with a color represent the ground truth pure pixels for a material. The ellipses with the same color represent the projected Gaussian components (twice the standard deviation along the major and minor axes, covering 86% of the total probability mass) for one endmember.

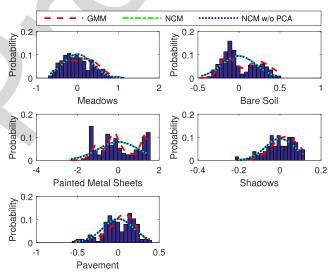


Fig. 8. Histograms of pure pixels for the Pavia University dataset and the estimated distributions from GMM and NCM when projected to 1 dimension.

941 samples in the ROI). Since painted metal sheets have multiple modes in the distribution, it should reflect a true difference between GMM and the other distributions. The abundances were sampled from a Dirichlet distribution so each pixel had random values. Also, an additive noise sampled from $\mathcal{N}(\mathbf{n}_n|\mathbf{0},\mathbf{D})$ was added to the mixed spectra, where the noise was assumed to be independent at different wavelengths, i.e. $\mathbf{D} = \operatorname{diag}\left(\sigma_1^2,\ldots,\sigma_B^2\right)$ while σ_k was again sampled from a uniform distribution on $[0,\sigma_Y]$.

We tested the algorithms for different σ_Y . The effects of priors were all removed in this case, i.e. $\beta_1 = 0$, $\beta_2 = 0$. Fig. 3 shows the box plots of abundance and endmember errors. We can see that GMM has small errors in general for different noise levels. NCM also has relatively small errors in most cases, but tends to produce large errors occasionally

696

698

700

702

704

706

707

708

711

712

713

715

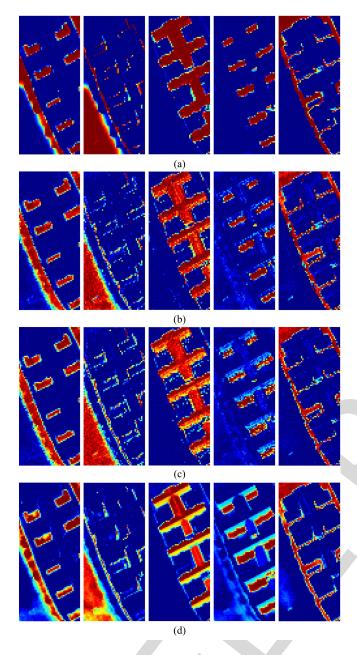


Fig. 9. Abundance maps for the Pavia University dataset. The corresponding endmembers from left to right are meadows, bare soil, painted metal sheets, shadows and pavement. (a) GMM. (b) NCM. (c) NCM w/o PCA. (d) BCM.

TABLE IV
ABUNDANCE AND ENDMEMBER ERRORS FOR PAVIA UNIVERSITY

$\times 10^{-4}$	GMM	NCM	NCM w/o PCA	BCM
Meadow	187 \ 44 ^a	405 \ 113	378 \ 114	711
Soil	175 \ 30	581 \ 68	507 \ 66	1049
Metal	476 \ 49	1236 \ 237	917 \ 349	1285
Shadow	44 \ 44	736 \ 48	914 \ 34	1287
Pavement	473 \ 39	1064 \ 114	333 \ 103	612
Mean	271 \ 41	804 \ 116	610 \ 133	989

^a the numbers in ".\." denote the abundance and endmember errors.

(4 out of 20 runs). NCM without PCA has very good results except for large noise, where it performed worst among all the methods. BCM has the largest errors overall. For the endmembers, although NCM or NCM without PCA sometimes has

691

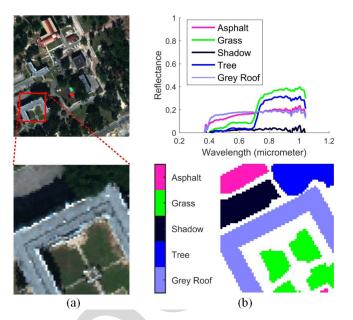


Fig. 10. (a) Original RGB image of the Mississippi Gulfport dataset with selected ROI and (b) Ground truth materials in the ROI with their mean spectra.

TABLE V $\\ Abundance \ and \ Endmember \ Errors \ for \ the \ Gulfport \ Dataset$

	$\times 10^{-4}$	GMM	NCM	NCM w/o PCA	BCM
ĺ	Asphalt	205 \ 52a	1693 \ 94	939 \ 59	1420
	Grass	169 \ 58	1982 \ 121	558 \ 65	2145
	Shadow	499 \ 49	1294 \ 68	921 \ 43	1315
Ì	Tree	1029 \ 89	2194 \ 234	1106 \ 185	2279
	Roof	908 \ 76	2143 \ 174	1234 \ 104	1657
	Mean	562 \ 65	1861 \ 138	952 \ 91	1763

^a the numbers in ".\." denote the abundance and endmember errors.

less errors than GMM, the difference is less than 0.005 hence negligible.

2) Unsupervised: We created two synthetic images in this case, the first was used to validate the ability to estimate the distribution parameters on scenes with regions of pure pixels, the second was used to validate the segmentation strategy on images with insufficient pure pixels. They were both of size 60×60 pixels and constructed from 4 endmember classes: limestone, basalt, concrete, asphalt, whose spectral signatures were highly differentiable. We assumed that the endmembers were sampled from GMMs following the example in Section II-C. The means of the GMMs were from the ASTER spectral library [50] (see Fig. 4(c) for their spectra) with slight constant changes, which determined a spectral range from $0.4~\mu m$ to $14~\mu m$, re-sampled into 200 values. The covariance matrices were constructed by $a_{jk}^2 \mathbf{I}_B + b_{jk}^2 \mathbf{u}_{jk} \mathbf{u}_{jk}^T$ where \mathbf{u}_{jk} was a unit vector controlling the major variation direction. For the first image, we assumed the 4 materials occupied the 4 quadrants of the square image as pure pixels. Then Gaussian smoothing was applied on each abundance map to make the boundary pixels of each quadrant be mixed by the neighboring materials. For the second image, we made the first material as background, the other materials randomly placed on this background. The procedure of generating the abundance

maps followed [37]: for each material (not as background), 150 Gaussian blobs were randomly placed, whose location and shape width were both sampled from Gaussian distributions. Finally, noise produced similar to above with $\sigma_Y = 0.001$ was added to the generated pixels. Fig. 4 shows the abundance maps, the original spectra of these materials, and the resulting color images by extracting the bands corresponding to wavelengths 488 nm, 556 nm, 693 nm.

The parameters of GMM were $r_{se} = 5$ for the two images, $\beta_1 = 0.1$, $\beta_2 = 0.1$ for the second image. Fig. 5 shows the histograms of ground truth pure pixels and the estimated distributions for the first image. The ground truth distribution is barely visible as most of the time it coincides with GMM. For limestone and asphalt, all the distributions are similar since the pure pixels are generated by a unimodal Gaussian. However, for basalt and concrete, GMM provides a more accurate estimation while the two NCMs seem inferior due to the single Gaussian assumption. The quantitative analysis in Table II implies a similar result by calculating the L_2 distance between the estimated distribution and the ground truth.

Table III shows the comparison of abundance errors from the two images. Since the second image is much more challenging than the first one, we can expect increased errors from all the methods. In general, the results of BCM and the two NCMs show slightly inferior abundances compared to GMM despite the fact that they have access to pure pixels in the image to train their models.

B. Pavia University

The Pavia University dataset was recorded by the Reflective Optics System Imaging Spectrometer (ROSIS) during a flight over Pavia, northern Italy. The dimension is 340 by 610 with a spatial resolution of 1.3 meters/pixel. It has 103 bands with wavelengths ranging from 430 nm to 860 nm. As Fig. 1 shows, the original image contains several man-made and natural materials. Considering that the whole dataset contains many different objects, we only performed experiments on the exemplar ROI (47 by 106) shown in Fig. 1, in which 5 endmembers, meadows, bare soil, painted metal sheets, shadows and pavement, are manually identified.

The parameter of GMM was $r_{se} = 2$. Fig. 6 shows the GMM in the wavelength-reflectance space, where we can see the centers and the major variations of the Gaussians. Fig. 7 shows the scatter plot of the results in the projected space. The scatter plot shows that the identified Gaussian components cover the ground truth pure pixels very well. For painted metal sheets, which has a broad range of pure pixels, it estimated 4 components to cover them. For shadows, only one component was estimated. Fig. 8 shows the histograms of pure pixels and the estimated distributions of GMM and NCMs. We can see that GMM matches the background histogram better than NCMs.

Fig. 9 shows the abundance map comparison. Comparing them with the ground truth shown in Fig. 1(a), we can see that BCM failed to estimate the pure pixels of painted metal sheets, although ground truth pure pixels were used for training.

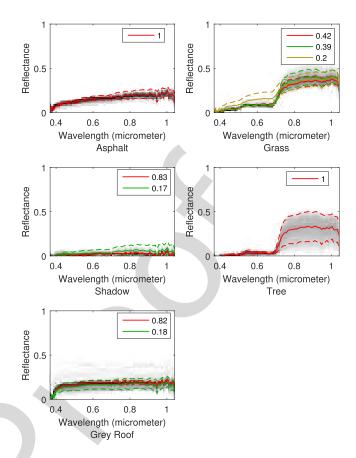


Fig. 11. Estimated GMM in the wavelength-reflectance space for the Mississippi Gulfport dataset. The background gray image and the curves have the same meaning as in Fig. 6.

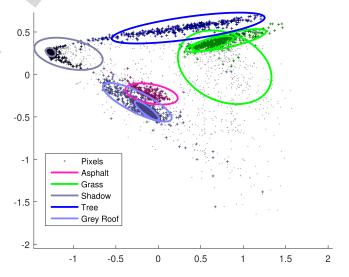


Fig. 12. Scatter plot of the Mississippi Gulfport dataset with the estimated GMM. The ellipses and the dots have the same meaning as in Fig. 7.

For example, the third and fourth abundance maps of BCM show that the pixels in the lower part of painted metal sheets are mixed with shadows, while the reduced reflectances are only caused by angle variation. The result of GMM not only shows sparse abundances for that region, but also interprets the boundary as a combination of neighboring materials. Since this

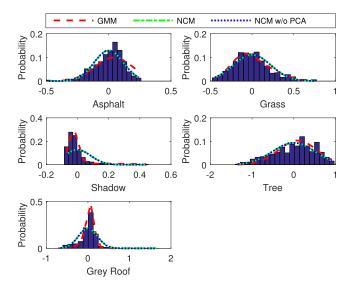


Fig. 13. Histograms of pure pixels for the Gulfport dataset and the estimated distributions from GMM and NCM when projected to 1 dimension.

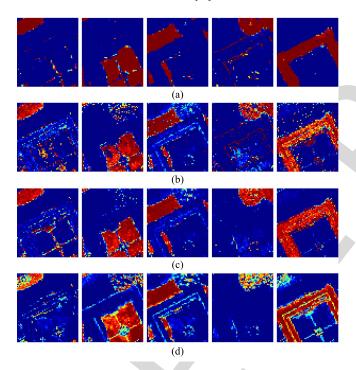


Fig. 14. Abundance maps for the Gulfport dataset. The corresponding endmembers from left to right are asphalt, grass, shadow, tree and grey roof. (a) GMM. (b) NCM. (c) NCM w/o PCA. (d) BCM.

dataset has a spatial spacing of 1.3 meters/pixel, we think this soft transition is more realistic than a simple segmentation. Although the results of NCMs look good in general, the abundances in a pure material region are inconsistent. The errors of abundances and endmembers for these algorithms are shown in Table IV, which implies that GMM performed best overall.

C. Mississippi Gulfport

The dataset was collected over the University of Southern Mississippis-Gulfpark Campus [51]. It is a 271 by 284 image with 72 bands corresponding to wavelengths 0.368 μm to 1.043 μm . The spatial resolution is 1 meter/pixel. The scene

contains several man-made and natural materials including sidewalks, roads, various types of building roofs, concrete, shrubs, trees, and grasses. Since the scene contains many cloths for target detection, we tried to avoid the cloths and selected a 58 by 65 ROI that contains 5 materials [52]. The original RGB image and the selected ROI are shown in Fig. 10(a) while the identified materials and the mean spectra are shown in (b).

The parameter of GMM was $r_{se} = 1$. Fig. 11 shows the GMM result in the wavelength-reflectance space and Fig. 12 shows the scatter plot. We can see that the estimated Gaussian components successfully cover the identified pure pixels. Fig. 13 shows the estimated distributions. Although there are no multiple peaks in any of the histograms, NCMs still do not fit the histograms of shadow and gray roof. In contrast, GMM gives a much better fit for these 2 endmember distributions.

Fig. 14 shows the abundance maps from different algorithms. We can see that GMM matches the ground truth in Fig. 10(b) best, followed by NCM without PCA. This is also verified in the quantitative analysis in Table V. Although NCM and BCM take ground truth pure pixels as input, the scattered dots for trees (fourth abundance map) in both of them and the incomplete region of grass for NCM (asphalt for BCM) show their insufficiency in this case.

V. DISCUSSION AND CONCLUSION

In this paper, we introduced a GMM approach to represent endmember variability, by observing that the identified pure pixels in real applications usually can not be well fitted by a unimodal distribution as in NCM or BCM. We solved several obstacles in linear unmixing using this distribution, including (i) deriving the conditional probability density function of the mixed pixel given each endmember modeled as GMM from two perspectives; (ii) estimating the abundances and endmember distributions by maximizing the log-likelihood with a prior enforcing abundance smoothness and sparsity; (iii) estimating the endmembers for each pixel given the abundances and distribution parameters. The results on synthetic and real datasets show superior accuracy compared to current popular methods like NCM, BCM. Here we have some final remarks.

A. Complexity

As analyzed in Section III-F, each iteration in the estimation of abundances has spatial complexity $O(|\mathcal{K}| NB^2)$ and time complexity $O(|\mathcal{K}| NB^3)$. For comparison, the implemented NCM has the same complexity but with $|\mathcal{K}|=1$. For the supervised synthetic dataset which contains 60 images, the total running time of GMM was 9709 seconds, on a desktop with a Intel Core i7-3820 CPU and 64 GB memory. For comparison, the running time of NCM, NCM without PCA, and BCM was 941, 50751, 62525 seconds respectively. In real applications, running GMM on the Pavia University and Mississippi Gulfport ROIs required 734 seconds and 97 seconds respectively for abundance estimation (24 seconds and 17 seconds for endmember estimation), compared to 40 and 34 seconds from NCM, 1389 and 396 seconds from

903

904

905

906

907

908

909 910

911

912

913

915

917

918

919

921

922

923

924

925

926

927

929

931

932

933

934

935

937

938

939

941

942

943

944

945

946

947

948

949

950

951

952

953

954

956

958

959

960

962

963

964

965

966

967

968

969

970

971

973

974

975

851

853

855

856

857

858

860

862

864

866

868

870

871

872

874

875

876

878

879

880

882

884

886

887

890

892

893

895

896

897

898

899

900

NCM without PCA, 1170 and 616 seconds from BCM. As analyzed, the main factors affecting the efficiency of GMM and NCMs are $|\mathcal{K}|$ and B.

B. Limitation

The complexity analysis leads to one limitation of the method. That is, the complexity grows exponentially with increasing numbers of components. This could cause problems for a large amount of pure pixels. To overcome this shortcoming, there are some empirical workarounds, such as reducing the number of components by introducing thresholds, or reducing the number of pure pixels to a fixed number by random sampling. Another limitation is that the proposed unsupervised version assumes presence of regions of pure pixels, which mostly happens in urban scenes. For scenes with a lot of mixed pixels, this assumption may not hold. Note that unsupervised unmixing is a very challenging problem. The previous works for this problem all assume several properties on the abundances and endmembers [21]-[23]. Hence, this limitation exists more or less in all the works on this problem. Finally, the method was only evaluated on real urban datasets with only ground truth on pure pixels: it is therefore unclear if the abundance estimation on mixed pixels is also accurate. This is due to lack of datasets and ground truth in the hyperspectral community. We plan to validate it on a more comprehensive dataset given in [31] in the future.

C. Future Work

The proposed GMM formulation has several applications that we can investigate in the future. First, in target detection, endmember variability may interfere with the target as well as the background [53]. By modeling the target or the background as spectra sampled from GMM distributions, we may devise more sophisticated and accurate target detection algorithms. Second, in fusion of hyperspectral and multispectral images, the LMM is usually used to overcome the underdetermined nature of the problem [54], [55]. However, the LMM does not hold in real scenarios as shown in this work. If we use the LMM with endmember variability, which is modeled by samples from GMM distributions, we may have a fusion algorithm that better fits the data. Finally, in estimating the noise or intrinsic dimension of hyperspectral images, simulated data are generated to quantify the results [46]. When these simulated data are created, usually the LMM is used without considering the endmember variability. Using the GMM formulation, we may generate distinct endmembers for each pixel and create more realistic synthetic data.

REFERENCES

- [1] M. Berman, H. Kiiveri, R. Lagerstrom, A. Ernst, R. Dunne, and J. F. Huntington, "ICE: A statistical approach to identifying endmembers in hyperspectral images," *IEEE Trans. Geosci. Remote Sens.*, vol. 42, no. 10, pp. 2085–2095, Oct. 2004.
- [2] J. M. P. Nascimento and J. M. Bioucas-Dias, "Vertex component analysis: A fast algorithm to unmix hyperspectral data," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 4, pp. 898–910, Apr. 2005.
- [3] A. Zare, P. Gader, O. Bchir, and H. Frigui, "Piecewise convex multiple-model endmember detection and spectral unmixing," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 5, pp. 2853–2862, May 2013.

- [4] J. M. Bioucas-Dias et al., "Hyperspectral unmixing overview: Geometrical, statistical, and sparse regression-based approaches," IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens., vol. 5, no. 2, pp. 354–379, Apr. 2012.
- [5] N. Keshava and J. F. Mustard, "Spectral unmixing," IEEE Signal Process. Mag., vol. 19, no. 1, pp. 44–57, Jan. 2002.
- [6] B. Hapke, "Bidirectional reflectance spectroscopy: 1. Theory," J. Geophys. Res., Solid Earth, vol. 86, no. B4, pp. 3039–3054, 1981.
- [7] A. Halimi, Y. Altmann, N. Dobigeon, and J.-Y. Tourneret, "Nonlinear unmixing of hyperspectral images using a generalized bilinear model," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 11, pp. 4153–4162, Nov. 2011.
- [8] B. Somers et al., "Nonlinear hyperspectral mixture analysis for tree cover estimates in orchards," *Remote Sens. Environ.*, vol. 113, no. 6, pp. 1183–1193, Feb. 2009.
- [9] R. Heylen and P. D. Gader, "Nonlinear spectral unmixing with a linear mixture of intimate mixtures model," *IEEE Geosci. Remote Sens. Lett.*, vol. 11, no. 7, pp. 1195–1199, Jul. 2014.
- [10] J. Broadwater and A. Banerjee, "A generalized kernel for areal and intimate mixtures," in *Proc. 2nd Workshop Hyperspectral Image Signal Process., Evol. Remote Sens. (WHISPERS)*, Jun. 2010, pp. 1–4.
- [11] J. Broadwater, R. Chellappa, A. Banerjee, and P. Burlina, "Kernel fully constrained least squares abundance estimates," in *Proc. IEEE Int. Geosci. Remote Sens. Symp. (IGARSS)*, Jul. 2007, pp. 4041–4044.
- [12] R. Heylen, M. Parente, and P. Gader, "A review of nonlinear hyperspectral unmixing methods," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 6, pp. 1844–1868, Jun. 2014.
- [13] B. Somers, G. P. Asner, L. Tits, and P. Coppin, "Endmember variability in spectral mixture analysis: A review," *Remote Sens. Environ.*, vol. 115, no. 7, pp. 1603–1616, 2011.
- [14] A. Zare and K. Ho, "Endmember variability in hyperspectral analysis: Addressing spectral variability during spectral unmixing," *IEEE Signal Process. Mag.*, vol. 31, no. 1, pp. 95–104, Jan. 2014.
- [15] X. Du, A. Zare, P. Gader, and D. Dranishnikov, "Spatial and spectral unmixing using the beta compositional model," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 6, pp. 1994–2003, Jun. 2014.
- [16] A. Zare and P. Gader, "PCE: Piecewise convex endmember detection," IEEE Trans. Geosci. Remote Sens., vol. 48, no. 6, pp. 2620–2632, Jun. 2010.
- [17] D. A. Roberts, M. Gardner, R. Church, S. Ustin, G. Scheer, and R. O. Green, "Mapping chaparral in the Santa Monica Mountains using multiple endmember spectral mixture models," *Remote Sens. Environ.*, vol. 65, no. 3, pp. 267–279, Sep. 1998.
- [18] A. Halimi, N. Dobigeon, and J.-Y. Tourneret, "Unsupervised unmixing of hyperspectral images accounting for endmember variability," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 4904–4917, Dec. 2015.
- [19] O. Eches, N. Dobigeon, C. Mailhes, and J.-Y. Tourneret, "Bayesian estimation of linear mixtures using the normal compositional model. Application to hyperspectral imagery," *IEEE Trans. Image Process.*, vol. 19, no. 6, pp. 1403–1413, Jun. 2010.
- [20] C. A. Bateson, G. P. Asner, and C. A. Wessman, "Endmember bundles: A new approach to incorporating endmember variability into spectral mixture analysis," *IEEE Trans. Geosci. Remote Sens.*, vol. 38, no. 2, pp. 1083–1094. Mar. 2000.
- [21] L. Drumetz, M.-A. Veganzones, S. Henrot, R. Phlypo, J. Chanussot, and C. Jutten, "Blind hyperspectral unmixing using an extended linear mixing model to address spectral variability," *IEEE Trans. Image Process.*, vol. 25, no. 8, pp. 3890–3905, Aug. 2016.
- [22] P.-A. Thouvenin, N. Dobigeon, and J.-Y. Tourneret, "Hyperspectral unmixing with spectral variability using a perturbed linear mixing model," *IEEE Trans. Signal Process.*, vol. 64, no. 2, pp. 525–538, Jan. 2016.
- [23] A. Halimi, P. Honeine, and J. M. Bioucas-Dias, "Hyperspectral unmixing in presence of endmember variability, nonlinearity, or mismodeling effects," *IEEE Trans. Image Process.*, vol. 25, no. 10, pp. 4565–4579, Oct. 2016.
- [24] B. Zhang, L. Zhuang, L. Gao, W. Luo, Q. Ran, and Q. Du, "PSO-EM: A hyperspectral unmixing algorithm based on normal compositional model," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 12, pp. 7782–7792, Dec. 2014.
- [25] O. Eches, N. Dobigeon, and J.-Y. Tourneret, "Estimating the number of endmembers in hyperspectral images using the normal compositional model and a hierarchical Bayesian algorithm," *IEEE J. Sel. Topics Signal Process.*, vol. 4, no. 3, pp. 582–591, Jun. 2010.

1056

1057

1059

1060

1061

1063

1064

1065

1066

1068

1069

1070

1071

1073

1074

1075

1077

1078

1079

1080

1082

1083

1085

1087

1089

1090

1092

1093

1094

1095

1097

1099

1100

1102

1104

1105

1106

1107

1108

1058 AO:3

- [27] J.-P. Combe et al., "Analysis of OMEGA/Mars express data hyperspectral data using a multiple-endmember linear spectral unmixing model (MELSUM): methodology and first results," Planetary Space Sci., vol. 56, no. 7, pp. 951–975, May 2008.
- [28] G. P. Asner and D. B. Lobell, "A biogeophysical approach for automated SWIR unmixing of soils and vegetation," Remote Sens. Environ., vol. 74, no. 1, pp. 99-112, Oct. 2000.
- [29] G. P. Asner and K. B. Heidebrecht, "Spectral unmixing of vegetation, soil and dry carbon cover in arid regions: Comparing multispectral and hyperspectral observations," Int. J. Remote Sens., vol. 23, no. 19, pp. 3939-3958, Oct. 2002.
- A. Castrodad, Z. Xing, J. B. Greer, E. Bosch, L. Carin, and G. Sapiro, "Learning discriminative sparse representations for modeling, source separation, and mapping of hyperspectral imagery," IEEE Trans. Geosci. Remote Sens., vol. 49, no. 11, pp. 4263-4281, Nov. 2011.
- E. B. Wetherley, D. A. Roberts, and J. P. McFadden, "Mapping spectrally similar urban materials at sub-pixel scales," Remote Sens. Environ., vol. 195, pp. 170-183, Jun. 2017.
- D. Stein, "Application of the normal compositional model to the analysis of hyperspectral imagery," in Proc. IEEE Workshop Adv. Techn. Anal. Remotely Sensed Data, Oct. 2003, pp. 44-51.
- [33] L. Tits, B. Somers, and P. Coppin, "The potential and limitations of a clustering approach for the improved efficiency of multiple endmember spectral mixture analysis in plant production system monitoring," IEEE Trans. Geosci. Remote Sens., vol. 50, no. 6, pp. 2273-2286, Jun. 2012.
- [34] M.-D. Iordache, L. Tits, J. M. Bioucas-Dias, A. Plaza, and B. Somers, "A dynamic unmixing framework for plant production system monitoring," IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens., vol. 7, no. 6, pp. 2016-2034, Jun. 2014.
- [35] L. Tits, B. Somers, W. Saeys, and P. Coppin, "Site-specific plant condition monitoring through hyperspectral alternating least squares unmixing," IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens., vol. 7, no. 8, pp. 3606-3618, Aug. 2014.
- [36] J. B. Lee, A. S. Woodyatt, and M. Berman, "Enhancement of high spectral resolution remote-sensing data by a noise-adjusted principal components transform," IEEE Trans. Geosci. Remote Sens., vol. 28, no. 3, pp. 295-304, May 1990.
- Y. Zhou, A. Rangarajan, and P. D. Gader, "A spatial compositional model for linear unmixing and endmember uncertainty estimation," IEEE Trans. Image Process., vol. 25, no. 12, pp. 5987–6002, Dec. 2016.
- [38] D. Achlioptas and F. McSherry, "On spectral learning of mixtures of distributions," in Learning Theory. Springer, 2005, pp. 458-469.
- N. Vlassis and A. Likas, "A greedy EM algorithm for Gaussian mixture learning," Neural Process. Lett., vol. 15, no. 1, pp. 77-87, Feb. 2002.
- [40] K. Lange, Optimization. Springer, 2013.
- X.-L. Meng and D. B. Rubin, "Maximum likelihood estimation via the ECM algorithm: A general framework," Biometrika, vol. 80, no. 2, pp. 267-278, 1993.
- D. P. Bertsekas, Nonlinear Programming. Belmont, MA, USA: Athena Scientific, 1999.
- C.-J. Lin, "Projected gradient methods for nonnegative matrix factorization," Neural Comput., vol. 19, no. 10, pp. 2756-2779, 2007.
- G. J. McLachlan and S. Rathnayake, "On the number of components in a Gaussian mixture model," Wiley Interdiscipl. Rev., Data Mining Knowl. Discovery, vol. 4, no. 5, pp. 341-355, 2014.
- [45] P. Smyth, "Model selection for probabilistic clustering using crossvalidated likelihood," Statist. Comput., vol. 10, no. 1, pp. 63-72,
- L. Gao, Q. Du, B. Zhang, W. Yang, and Y. Wu, "A comparative study on linear regression-based noise estimation for hyperspectral imagery, IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens., vol. 6, no. 2, pp. 488-498, Apr. 2013.
- [47] R. E. Roger, "Principal Components transform with simple, automatic noise adjustment," Int. J. Remote Sens., vol. 17, no. 14, pp. 2719-2727,
- N. Guan, D. Tao, Z. Luo, and B. Yuan, "Manifold regularized discriminative nonnegative matrix factorization with fast gradient descent," IEEE Trans. Image Process., vol. 20, no. 7, pp. 2030-2048, 1049 Jul. 2011.
 - C. Bishop, Pattern Recognition and Machine Learning. New York, NY, USA: Springer, 2006

- [50] A. M. Baldridge, S. J. Hook, C. I. Grove, and G. Rivera, "The ASTER spectral library version 2.0," Remote Sens. Environ., vol. 113, no. 4, pp. 711-715, 2009.
- [51] P. Gader, A. Zare, R. Close, J. Aitken, and G. Tuell, "MUUFL Gulfport hyperspectral and LiDAR airborne data set," Univ. Florida, Gainesville, FL, USA, Tech. Rep. REP-2013-570, 2013.
- [52] X. Du and A. Zare, "Technical report: Scene label ground truth map for MUUFL Gulfport data set," Univ. Florida, Gainesville, FL, USA, Tech. Rep. 20170417, 2017.
- [53] C. Jiao and A. Zare, "Functions of multiple instances for learning target signatures," IEEE Trans. Geosci. Remote Sens., vol. 53, no. 8, pp. 4670–4686, Aug. 2015.
- [54] N. Yokoya, T. Yairi, and A. Iwasaki, "Coupled nonnegative matrix factorization unmixing for hyperspectral and multispectral data fusion," IEEE Trans. Geosci. Remote Sens., vol. 50, no. 2, pp. 528-537, Feb. 2012.
- Q. Wei, N. Dobigeon, and J.-Y. Tourneret, "Fast fusion of multi-band images based on solving a Sylvester equation," IEEE Trans. Image Process., vol. 24, no. 11, pp. 4109-4121, Nov. 2015.



Yuan Zhou received the B.E degree in software engineering and the M.E. degree in computer application technology from the Huazhong University of Science and Technology, Wuhan, China, in 2008 and 2011, respectively. Since 2013, he is currently pursuing the Ph.D. degree with the Department of CISE, University of Florida, Gainesville, FL, USA. He was with Shanghai UIH as a Software Engineer for two years. His research interests include image processing, computer vision, and machine learning.



Anand Rangarajan is currently with the Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL, USA. His research interests are machine learning, computer vision, and the scientific study of consciousness.



Paul D. Gader (M'86–SM'09–F'11) received the Ph.D. degree in mathematics for image-processingrelated research from the University of Florida, Gainesville, FL, USA, in 1986. He was a Senior Research Scientist with Honeywell, a Research Engineer and a Manager with the Environmental Research Institute of Michigan, Ann Arbor, MI, USA, and a Faculty Member with the University of Wisconsin, Oshkosh, WI, USA, the University of Missouri, Columbia, MO, USA, and the University of Florida, FL, USA, where he is currently a

Professor of Computer and Information Science and Engineering. He was a Summer Student Fellow at the Eglin Air Force Base involving in algorithms for the detection of bridges in forward-looking infrared imagery, where he performed his first research in image processing in 1984. He has been involved in a wide variety of theoretical and applied research problems including fast computing with linear algebra, mathematical morphology, fuzzy sets, Bayesian methods, handwriting recognition, automatic target recognition, biomedical image analysis, landmine detection, human geography, and hyperspectral and light detection, and ranging image analysis projects. He has authored or co-authored hundreds of refereed journal and conference papers.

977

978

979

980 981

982 983

984

985

986

987

988

989

990

991

992

993

994

995

996 997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024 1025

1026

1027

1028

1029

1030 1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1050

1051

1052

AO:2

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES

PLEASE NOTE: We cannot accept new source files as corrections for your paper. If possible, please annotate the PDF proof we have sent you with your corrections and upload it via the Author Gateway. Alternatively, you may send us your corrections in list format. You may also upload revised graphics via the Author Gateway.

AQ:1 = Please note that there were discrepancies between the accepted pdf [gmm_var_journal_ieee_final.pdf] and the [gmm_var_journal_ieee_final.tex] in the sentence on lines 188–195, 623–632, and references. We have followed [gmm_var_journal_ieee_final.tex].

AQ:2 = Please provide the publisher location for refs. [38] and 140].

AQ:3 = Please provide the department name for refs. [51] and [52].





A Gaussian Mixture Model Representation of Endmember Variability in Hyperspectral Unmixing

Yuan Zhou[®], Student Member, IEEE, Anand Rangarajan, Member, IEEE, and Paul D. Gader, Fellow, IEEE

Abstract—Hyperspectral unmixing while considering endmember variability is usually performed by the normal compositional model, where the endmembers for each pixel are assumed to be sampled from unimodal Gaussian distributions. However, in real applications, the distribution of a material is often not Gaussian. In this paper, we use Gaussian mixture models (GMM) to represent endmember variability. We show, given the GMM starting premise, that the distribution of the mixed pixel (under the linear mixing model) is also a GMM (and this is shown from two perspectives). The first perspective originates from random variable transformations and gives a conditional density function of the pixels given the abundances and GMM parameters. With proper smoothness and sparsity prior constraints on the abundances, the conditional density function leads to a standard maximum a posteriori (MAP) problem which can be solved using generalized expectation maximization. The second perspective originates from marginalizing over the endmembers in the GMM, which provides us with a foundation to solve for the endmembers at each pixel. Hence, compared to the other distribution based methods, our model can not only estimate the abundances and distribution parameters, but also the distinct endmember set for each pixel. We tested the proposed GMM on several synthetic and real datasets, and showed its potential by comparing it to current popular methods.

12

13

14

17

18

21

22

23

27

28

29

30

31

32

33

Index Terms—Endmember extraction, endmember variability, hyperspectral image analysis, linear unmixing, Gaussian mixture model.

I. INTRODUCTION

THE formation of hyperspectral images can be simplified by the *linear mixing model* (LMM), which assumes that the physical region corresponding to a pixel contains several pure materials, so that each material contributes a fraction of its spectra based on area to the final spectra of the pixel. Hence, the observed spectra $\mathbf{y}_n \in \mathbb{R}^B$, $n = 1, \ldots, N$ (B is the number of wavelengths and N is the number of pixels) is a (non-negative) linear combination of the pure material

Manuscript received June 30, 2017; revised November 23, 2017; accepted January 10, 2018. This work was supported by NSF IIS under Grant 1743050. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Jocelyn Chanussot. (Corresponding author: Yuan Zhou.)

The authors are with the Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL 32611 USA (e-mail: yuan@cise.ufl.edu; anand@cise.ufl.edu; pgader@cise.ufl.edu).

This paper has supplementary downloadable material available a http://ieeexplore.ieee.org., provided by the author.

Color versions of one or more of the figures in this paper are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/TIP.2018.2795744

(called *endmember*) spectra $\mathbf{m}_j \in \mathbb{R}^B$, j = 1, ..., M (M is the number of endmembers), i.e.

$$\mathbf{y}_{n} = \sum_{j=1}^{M} \mathbf{m}_{j} \alpha_{nj} + \mathbf{n}_{n}, \text{ s.t. } \alpha_{nj} \ge 0, \quad \sum_{j=1}^{M} \alpha_{nj} = 1, \quad (1)$$

where a_{nj} is the proportion (called *abundance*) for the *j*th endmember at the *n*th pixel (with the positivity and sumto-one constraint) and $\mathbf{n}_n \in \mathbb{R}^B$ is additive noise. Here, the endmember set $\{\mathbf{m}_j : j=1,\ldots,M\}$ is fixed for all the pixels. This model simplifies the unmixing problem to a matrix factorization one, leading to efficient computation and simple algorithms such as iterative constrained endmembers (ICE), vertex component analysis (VCA), piecewise convex multiplemodel endmember detection (PCOMMEND) [1]–[3] etc., which receive comprehensive reviews in [4] and [5].

However, in practice the LMM may not be valid in many real scenarios. Even for a *pure* pixel that only contains one material, its spectrum may not be consistent over the whole image. This is due to several factors such as atmospheric conditions, topography and intrinsic variability. For example, in vegetation, multiple scattering and biotic variation (e.g. differences in biochemistry and water content) cause different reflectances among the same species. For urban scenes, the incidence and emergence angles could be different for the same roof, causing different reflectances. For minerals, the spectroscopy model developed by Hapke also considers the porosity and roughness of the material as variable [6].

50

55

61

65

73

74

76

In the first and third example above, Eq. (1) can be generalized to a more abstract form $\mathbf{y}_n = F(\{\mathbf{m}_j, \alpha_{nj} : j=1,\ldots M\})$, which leads to nonlinear mixing models. For example, Halimi et al. [7] used bilinear models to handle the vegetation case, which was also investigated using several different nonlinear functions [8]. In [9], the Hapke model was used to model intimate interaction among minerals. There are also works that use kernels for flexible nonlinear mixing [10], [11]. A panoply of nonlinear models can be found in the review article [12]. We note that in these models, a fixed endmember set is still assumed while using a more complicated unmixing model.

While nonlinear models abound lately, it is still difficult to account for all the scenarios. On the contrary, the LMM still has physical significance with the intuitive area assumption. To model real scenarios more accurately, researchers have

87

93

97

100

101

102

104

106

108

110

112

115

116

117

119

120

121

123

125

126

127

129

130

taken another route by generalizing Eq. (1) to

$$\mathbf{y}_n = \sum_{j=1}^M \mathbf{m}_{nj} \, \alpha_{nj} + \mathbf{n}_n, \tag{2}$$

where $\{\mathbf{m}_{nj} \in \mathbb{R}^B : j = 1, ..., M\}, n = 1, ..., N$ could be different for each n, i.e. the endmember spectra for each pixel could be different. This is called endmember variability, and has also received a lot of attention in the community [13], [14]. Note that given $\{\mathbf{y}_n\}$, inferring $\{\mathbf{m}_{nj}, \alpha_{nj}\}$ is a much more difficult problem than inferring $\{\mathbf{m}_j, \alpha_{nj}\}$ in Eq. (1). Hence, in many papers $\{\mathbf{m}_{nj}\}$ are assumed to be from a spectral library, which is usually called *supervised unmixing* [15]–[17]. On the other hand, if the endmember spectra are to be extracted from the image, we call them unsupervised unmixing models [18]–[20]. Obviously, unsupervised unmixing more challenging than its supervised counterpart and hence more assumptions are used in this case, such as the spatial smoothness of abundances and endmember variability [21]-[23], small mutual distance between the endmembers [22], small magnitude or spectral smoothness of the endmember variability [22], [23].

We can also categorize the papers on endmember variability by how this variability is modeled. In the review paper [14], it can be modeled as a endmember set [17], [20] or as a distribution [24]–[26]. One of the widely used set based methods is multiple endmember spectral mixture analysis (MESMA) [17], which tries every endmember combination and selects the one with the smallest error. There are many variations to the original MESMA. For example, the multiple-endmember linear spectral unmixing model (MELSUM) solves the linear equations directly using the pseudo-inverse and discards the solutions with negative abundances [27]; automatic Monte Carlo unmixing (AutoMCU) picks random combinations for unmixing and averages the resulting abundances as the final results [28], [29]. Besides MESMA variants, there are also many other set based methods. For example, endmember bundles form bundles from automated extracted endmembers, take minimum and maximum abundances from bundle based unmixing, and average them as final abundances [20]; sparse unmixing imposes a sparsity constraint on the abundances based on endmembers composed of all spectra from the spectral library [30]. A comprehensive review can be found in [13] and [14]. One disadvantage of set based methods is that their complexity increases exponentially with increasing library size hence in practice a laborious library reduction approach may be required [31].

The distribution based approaches assume that the endmembers for each pixel are sampled from probability distributions [e.g. Gaussian, a.k.a. normal compositional model (NCM)], and hence embrace large libraries while being numerically tractable [15], [32]. Here, we give an overview of NCM because of its simplicity and popularity [16], [18], [19]. Suppose the *j*th endmember at the *n*th pixel follows a Gaussian distribution $p(\mathbf{m}_{nj}) = \mathcal{N}(\mathbf{m}_{nj}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ where $\boldsymbol{\mu}_j \in \mathbb{R}^B$ and $\boldsymbol{\Sigma}_j \in \mathbb{R}^{B \times B}$, and the additive noise also follows a Gaussian distribution $p(\mathbf{n}_n) = \mathcal{N}(\mathbf{n}_n|\mathbf{0},\mathbf{D})$ where \mathbf{D} is the noise covariance matrix. The random variable transformation (r.v.t.)

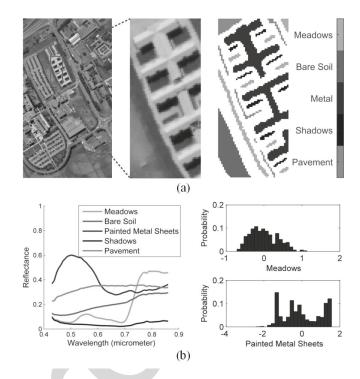


Fig. 1. (a) Original Pavia University image and selected ROI with its ground truth image. (b) Mean spectra of the identified 5 endmembers and histograms of meadows and painted metal sheets (shadow is termed as endmember to conform with the LMM though the area under shadow can be any material). PCA is used to project the multidimensional pixels to single values which are counted in the histograms. Although the histogram of meadows may appear to be a Gaussian distribution, that of painted metal sheets is obviously neither a unimodal Gaussian or Beta distribution.

(2) suggests that the probability density function of \mathbf{y}_n can be derived as

$$p(\mathbf{y}_n|\boldsymbol{\alpha}_n,\boldsymbol{\Theta},\mathbf{D}) = \mathcal{N}\left(\mathbf{y}_n|\sum_{j=1}^{M} \alpha_{nj}\boldsymbol{\mu}_j, \sum_{j=1}^{M} \alpha_{nj}^2 \boldsymbol{\Sigma}_j + \mathbf{D}\right), \quad (3)$$

137

139

141

142

143

144

145

146

147

148

150

151

152

154

where $\alpha_n := [\alpha_{n1}, \dots, \alpha_{nM}]^T$, $\Theta := \{\mu_j, \Sigma_j : j = 1, \dots, M\}$. The conditional density function in (3) is usually embedded in a Bayesian framework such that we can incorporate priors and also estimate hyperparameters. Then, NCM uses different optimization approaches, e.g. expectation maximization [32], sampling methods [18], [19], [25], particle swarm optimization [24], to determine the parameters $\{\mu_i, \Sigma_i\}$ and $\{\alpha_{ni}\}$.

There are few papers that use other distributions. In [15], X. Du *et al.* note that the Gaussian distribution may allow negative values which are not realistic. In addition, the real distribution may be skewed. Hence, they introduce a Beta compositional model (BCM) to model the variability. The problem is that the true distribution may not be well approximated by any unimodal distribution. Consider the Pavia University dataset shown in Fig. 1, where the multidimensional pixels are projected to one dimension to afford better visualization. Among the manually identified materials, we can see that although the histogram of meadows may look like a Gaussian distribution, that of painted metal sheets has multiple peaks and cannot be approximated by either a Gaussian or

205

207

211

212

213

215

220

221

223

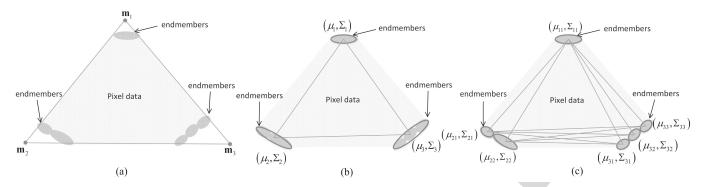


Fig. 2. Comparison of the mechanisms among LMM, NCM and GMM. We have 3 endmembers represented by the darken gray areas. LMM tries to find a set of endmembers that fit the pixel data. NCM tries to find a set of Gaussian centers that fit the pixel data, with error weighted by the covariance matrices. GMM tries to find Gaussian centers such that all their linear combinations fit the pixel data, with each weighted by the prior π_k . We may use 6 endmembers with NCM, but then the prior information is lost.

Beta distribution. This is due to different angles of these sheets on the roof. Since each piece of metal sheet is tilted, it forms a cluster of reflectances which contributes to a peak in the histogram. This example shows that we should use a more flexible distribution to represent the endmember variability.

158

159

160

161

162

163

164

165

167

169

171

173

174

175

176

178

180

181

182

183

184

185

186

188

189

190

191

192

193

194

195

197

AO:1

In this paper, we use a mixture of Gaussians to approximate any distribution that an endmember may exhibit, and solve the LMM by considering endmember variability. In a nutshell, the Gaussian mixture model (GMM) models $p(\mathbf{m}_{ni})$ by a mixture of Gaussians, say $p(\mathbf{m}_{ni}) =$ $\sum_{k} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right)$, and then obtains the distribution of y_n by the r.v.t. (2), which turns out to be another mixture of Gaussians and can be used for inference of the unknown parameters. Here, we briefly explain how GMM works intuitively by comparing it to the NCM with the details given later. The maximum likelihood estimate (MLE) of NCM (using (3)) aims to find $\{\mu_i\}$ such that its linear combination matches y_n . Contrary to NCM, GMM aims to find $\{\mu_{jk}\}$ such that all of its linear combinations match y_n . Suppose we have μ_{11} , μ_{21} , μ_{22} , μ_{31} , μ_{32} , μ_{33} : then there are 6 combinations as explained in Fig. 2, but with emphasis weighted by $\{\pi_{ik}\}$ which determines the prior probability of each linear combination.

Based on the GMM formulation, we propose a supervised version and an unsupervised version for unmixing. The supervised version takes a library as input and estimates the abundances. The unsupervised version assumes that there are regions of pure pixels, hence segments the image first to get pure pixels and then performs unmixing. Another advantage over the other distribution based methods is that we can also estimate the endmembers for each pixel, which is not achievable by NCM or BCM. Note that estimating endmembers for each pixel is generally common in non-distribution methods, both from the signal processing community [21]–[23] or the remote sensing community [17], [27]. But it is often achieved in the context of least-squares based unmixing [33]–[35], unlike what we propose here using distribution based unmixing.

Notation: As usual, $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the multivariate Gaussian density function with center $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a matrix with m rows and n columns. The Hadamard product of two matrices (elementwise multiplication) is denoted by \circ while the Kronecker

product is denoted by \otimes . (A)_{jk} denotes the element at the jth row and kth column of matrix **A**. (**A**)_j denotes the jth row of **A** transposed (treating **A** as a vector), i.e. for **A** = $[\mathbf{a}_1, \dots \mathbf{a}_n]^T$, (**A**)_j = \mathbf{a}_j . vec (**A**) denotes the vectorization of **A**, i.e. concatenating the columns of **A**. $\delta_{jk} = 1$ when j = k and 0 otherwise. $\mathbb{E}_{\mathbf{x}}(f(\mathbf{x}))$ is the expected value of $f(\mathbf{x})$ given random variable **x**. We use $i = \sqrt{-1}$ instead of as an index throughout the paper.

II. MATHEMATICAL PRELIMINARIES

A. Linear Combination of GMM Random Variables

To use the Gaussian mixture model to model endmember variability, we start by assuming that \mathbf{m}_{nj} follows a Gaussian mixture model (GMM) and the noise also follows a Gaussian distribution. The distribution of \mathbf{y}_n is obtained using the following theorem.

Theorem 1: If the random variable \mathbf{m}_{nj} has a density function

$$p\left(\mathbf{m}_{nj}|\boldsymbol{\Theta}\right) := f_{\mathbf{m}_{j}}\left(\mathbf{m}_{nj}\right) = \sum_{k=1}^{K_{j}} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj}|\boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right), \quad (4) \quad {}_{21}$$

s.t. $\pi_{jk} \geq 0$, $\sum_{k=1}^{K_j} \pi_{jk} = 1$, with K_j being the number of components, π_{jk} ($\mu_{jk} \in \mathbb{R}^B$ or $\Sigma_{jk} \in \mathbb{R}^{B \times B}$) being the weight (mean or covariance matrix) of its kth Gaussian component, $\Theta := \{\pi_{jk}, \mu_{jk}, \Sigma_{jk} : j = 1, \dots, M, k = 1, \dots, K_j\}$, $\{\mathbf{m}_{nj} : j = 1, \dots, M\}$ are independent, and the random variable \mathbf{n}_n has a density function $p(\mathbf{n}_n) := \mathcal{N}(\mathbf{n}_n | \mathbf{0}, \mathbf{D})$, then the density function of \mathbf{y}_n given by the r.v.t. $\mathbf{y}_n = \sum_{j=1}^M \mathbf{m}_{nj} \alpha_{nj} + \mathbf{n}_n$ is another GMM

$$p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right) = \sum_{\mathbf{k}\in\mathcal{K}} \pi_{\mathbf{k}} \mathcal{N}\left(\mathbf{y}_{n}|\boldsymbol{\mu}_{n\mathbf{k}},\boldsymbol{\Sigma}_{n\mathbf{k}}\right),\tag{5}$$

where $K := \{1, ..., K_1\} \times \{1, ..., K_2\} \times \cdots \times \{1, ..., K_M\}$ is the Cartesian product of the M index sets, $\mathbf{k} := (k_1, ..., k_M) \in \mathcal{K}$, $\pi_{\mathbf{k}} \in \mathbb{R}$, $\mu_{n\mathbf{k}} \in \mathbb{R}^B$, $\Sigma_{n\mathbf{k}} \in \mathbb{R}^{B \times B}$ are defined by

$$\pi_{\mathbf{k}} := \prod_{j=1}^{M} \pi_{jk_j}, \quad \mu_{n\mathbf{k}} := \sum_{j=1}^{M} \alpha_{nj} \mu_{jk_j},$$
 229

$$\mathbf{\Sigma}_{n\mathbf{k}} := \sum_{j=1}^{M} \alpha_{nj}^2 \mathbf{\Sigma}_{jk_j} + \mathbf{D}. \tag{6}$$

The proof is detailed using a characteristic function (c.f.) approach.

We first consider the distribution of the intermediate variable $\mathbf{z}_n = \sum_{j=1}^{M} \mathbf{m}_{nj} \alpha_{nj}$. The c.f. of $f_{\mathbf{m}_j}$ in (4), $\phi_{\mathbf{m}_j}$ (t) : $\mathbb{R}^B \to \mathbb{C}$, is given by

$$\phi_{\mathbf{m}_{j}}(\mathbf{t}) = \mathbb{E}_{\mathbf{m}_{j}}\left(e^{i\mathbf{t}^{T}\mathbf{x}}\right) = \int_{\mathbb{R}^{B}} e^{i\mathbf{t}^{T}\mathbf{x}} f_{\mathbf{m}_{j}}(\mathbf{x}) d\mathbf{x}$$

$$= \sum_{k=1}^{K_{j}} \pi_{jk} \int_{\mathbb{R}^{B}} e^{i\mathbf{t}^{T}\mathbf{x}} \mathcal{N}\left(\mathbf{x} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right) d\mathbf{x}$$

$$= \sum_{k=1}^{K_{j}} \pi_{jk} \phi_{jk}(\mathbf{t}), \tag{7}$$

where ϕ_{jk} (t) denotes the c.f. of the Gaussian distribution $\mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}_{jk},\boldsymbol{\Sigma}_{jk}\right)$ as

$$\phi_{jk}(\mathbf{t}) := \exp\left(i\mathbf{t}^T \boldsymbol{\mu}_{jk} - \frac{1}{2}\mathbf{t}^T \boldsymbol{\Sigma}_{jk}\mathbf{t}\right).$$
 (8)

Assuming $\mathbf{m}_{n1}, \dots, \mathbf{m}_{nM}$ are independent, we can obtain the c.f. of the linear combination of these \mathbf{m}_{nj} by multiplying (7) as

 $\phi_{\mathbf{z}_n}\left(\mathbf{t}
ight)$

$$= \phi_{\mathbf{m}_{n1}\alpha_{n1} + \dots + \mathbf{m}_{nM}\alpha_{nM}} (\mathbf{t}) = \prod_{j=1}^{M} \phi_{\mathbf{m}_{j}} (\alpha_{nj} \mathbf{t})$$

$$=\sum_{k_1=1}^{K_1}\cdots\sum_{k_M=1}^{K_M}\pi_{1k_1}\cdots\pi_{Mk_M}\phi_{1k_1}(\alpha_{n1}\mathbf{t})\cdots\phi_{Mk_M}(\alpha_{nM}\mathbf{t}).$$

Let K, K, π_k be defined as in Theorem 1. We can write the above multiple summations in an elegant way:

$$\phi_{\mathbf{z}_n}(\mathbf{t}) = \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \phi_{n\mathbf{k}}(\mathbf{t}), \tag{9}$$

where $\pi_{\mathbf{k}} \geq 0$, $\sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} = 1$ and

 $\phi_{n\mathbf{k}}\left(\mathbf{t}\right):=\phi_{1k_{1}}\left(\alpha_{n1}\mathbf{t}\right)\cdots\phi_{Mk_{M}}\left(\alpha_{nM}\mathbf{t}\right)$

$$= \exp \left\{ i \mathbf{t}^T \left(\sum_{j=1}^M \alpha_{nj} \boldsymbol{\mu}_{jk_j} \right) - \frac{1}{2} \mathbf{t}^T \left(\sum_{j=1}^M \alpha_{nj}^2 \boldsymbol{\Sigma}_{jk_j} \right) \mathbf{t} \right\},\,$$

where (8) is used. Since $\phi_{n\mathbf{k}}(\mathbf{t})$ also has a form of c.f. of a Gaussian distribution, the corresponding distribution turns out to be $\mathcal{N}\left(\mathbf{x}|\sum_{j}\alpha_{nj}\boldsymbol{\mu}_{jk_{j}},\sum_{j}\alpha_{nj}^{2}\boldsymbol{\Sigma}_{jk_{j}}\right)$. Hence, the distribution of \mathbf{z}_{n} can be obtained by the Fourier transform of (9)

$$f_{\mathbf{z}_{n}}(\mathbf{z}_{n}) = \frac{1}{(2\pi)^{B}} \int_{\mathbb{R}^{B}} e^{-i\mathbf{t}^{T}\mathbf{z}_{n}} \phi_{\mathbf{z}_{n}}(\mathbf{t}) d\mathbf{t}$$

$$= \frac{1}{(2\pi)^{B}} \int_{\mathbb{R}^{B}} e^{-i\mathbf{t}^{T}\mathbf{z}_{n}} \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \phi_{n\mathbf{k}}(\mathbf{t}) d\mathbf{t}$$

$$= \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \mathcal{N} \left(\mathbf{z}_{n} | \sum_{i=1}^{M} \alpha_{nj} \boldsymbol{\mu}_{jk_{j}}, \sum_{i=1}^{M} \alpha_{nj}^{2} \boldsymbol{\Sigma}_{jk_{j}} \right), \quad (10)$$

which is still a mixture of Gaussians.

After finding the distribution of the linear combination, we can add the noise term to find the distribution of y_n . Suppose the noise also follows a Gaussian distribution,

 $p(\mathbf{n}_n) := f_{\mathbf{n}_n}(\mathbf{n}_n) = \mathcal{N}(\mathbf{n}_n|\mathbf{0},\mathbf{D})$, where \mathbf{D} is the noise covariance matrix. We assume that the noise at different wavelengths is independent (σ_k^2) being the noise variance of the kth band), i.e. $\mathbf{D} = \mathrm{diag}(\sigma_1^2,\sigma_2^2,\ldots,\sigma_B^2) \in \mathbb{R}^{B\times B}$ (if it is not independent, the noise can actually be easily whitened to be independent as in [36]). Its c.f. has the following form

$$\phi_{\mathbf{n}_n}(\mathbf{t}) = \exp\left(-\frac{1}{2}\mathbf{t}^T\mathbf{D}\mathbf{t}\right) \tag{11}$$

by (8). Then the c.f. of \mathbf{y}_n can be obtained by multiplying (9) and (11) (as \mathbf{z}_n and \mathbf{n}_n are independent)

$$\phi_{\mathbf{y}_{n}}\left(\mathbf{t}\right)=\phi_{\mathbf{z}_{n}}\left(\mathbf{t}\right)\phi_{\mathbf{n}_{n}}\left(\mathbf{t}\right)=\sum_{\mathbf{k}\in\mathcal{K}}\pi_{\mathbf{k}}\phi_{\mathbf{n}_{n}}\left(\mathbf{t}\right)\phi_{n\mathbf{k}}\left(\mathbf{t}\right)$$
 274

$$= \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \exp \left\{ i \mathbf{t}^T \boldsymbol{\mu}_{n\mathbf{k}} - \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma}_{n\mathbf{k}} \mathbf{t} \right\},$$
 275

where $\mu_{n\mathbf{k}}$ and $\Sigma_{n\mathbf{k}}$ are defined in (6). Finally, the distribution of \mathbf{y} can be shown to be (5) by the Fourier transform again as in (10).

If $K = \{1\} \times \{1\} \times \dots \times \{1\}$, i.e. each endmember has only one Gaussian component, we have $\pi_{11} = 1, \dots, \pi_{M1} = 1$, then $\pi_{\mathbf{k}} = \pi_{11} \cdots \pi_{M1} = 1$. The distribution of \mathbf{y}_n becomes

$$p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right) = \mathcal{N}\left(\mathbf{y}_{n}|\sum_{j=1}^{M}\alpha_{nj}\boldsymbol{\mu}_{j1},\sum_{j=1}^{M}\alpha_{nj}^{2}\boldsymbol{\Sigma}_{j1} + \mathbf{D}\right),$$
(12) 283

which is exactly the NCM in (3).

B. Another Perspective

Theorem 1 obtains the density of each pixel by directly performing a r.v.t. based on the LMM, which can be used to estimate the abundances and distribution parameters. Here, we will obtain the density from another perspective, which provides a foundation to estimate the endmembers for each pixel. Again, let the noise follow the density function $p(\mathbf{n}_n) := \mathcal{N}(\mathbf{n}_n | \mathbf{0}, \mathbf{D})$. Considering $\{\mathbf{m}_{nj}\}$ and $\{\alpha_{nj}\}$ as fixed values, the r.v.t. $\mathbf{y}_n = \sum_j \mathbf{m}_{nj} \alpha_{nj} + \mathbf{n}_n$ implies that the density of \mathbf{y}_n is given by

$$p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\mathbf{M}_{n},\mathbf{D}\right) = \mathcal{N}\left(\mathbf{y}_{n}|\sum_{j}\mathbf{m}_{nj}\alpha_{nj},\mathbf{D}\right)$$
(13) 295

where $\mathbf{M}_n = [\mathbf{m}_{n1}, \dots, \mathbf{m}_{nM}]^T \in \mathbb{R}^{M \times B}$ are the endmembers for the *n*th pixel. We have the following theorem which gives the same result as in Theorem 1.

Theorem 2: If the random variables $\{\mathbf{m}_{nj}: j=1,\ldots,M\}$ follow GMM distributions

$$p\left(\mathbf{m}_{nj}|\mathbf{\Theta}\right) := \sum_{k=1}^{K_j} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj}|\boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right),$$
 301

and they are independent, i.e.

$$p\left(\mathbf{M}_{n}|\mathbf{\Theta}\right) = \prod_{j=1}^{M} p\left(\mathbf{m}_{nj}|\mathbf{\Theta}\right),\tag{14}$$

350

352

355

357

364

366

371

372

376

378

380

381

TABLE I
VALUES FOR THE VARIOUS QUANTITIES IN THE SIMPLE EXAMPLE

k	$\pi_{\mathbf{k}}$	$\mu_{n\mathbf{k}}$ in (6)
(1,1,1,1)	0.06	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{21} + \alpha_{n3}\boldsymbol{\mu}_{31} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1, 2, 1, 1)	0.14	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{22} + \alpha_{n3}\boldsymbol{\mu}_{31} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1,1,2,1)	0.12	$\alpha_{n1}\mu_{11} + \alpha_{n2}\mu_{21} + \alpha_{n3}\mu_{32} + \alpha_{n4}\mu_{41}$
(1, 2, 2, 1)	0.28	$\alpha_{n1}\mu_{11} + \alpha_{n2}\mu_{22} + \alpha_{n3}\mu_{32} + \alpha_{n4}\mu_{41}$
(1,1,3,1)	0.12	$\alpha_{n1}\boldsymbol{\mu}_{11} + \alpha_{n2}\boldsymbol{\mu}_{21} + \alpha_{n3}\boldsymbol{\mu}_{33} + \alpha_{n4}\boldsymbol{\mu}_{41}$
(1,2,3,1)	0.28	$\alpha_{n1}\mu_{11} + \alpha_{n2}\mu_{22} + \alpha_{n3}\mu_{33} + \alpha_{n4}\mu_{41}$

then the conditional density $p(\mathbf{y}_n|\boldsymbol{\alpha}_n, \boldsymbol{\Theta}, \mathbf{D})$ obtained by marginalizing \mathbf{M}_n in $p(\mathbf{y}_n, \mathbf{M}_n|\boldsymbol{\alpha}_n, \boldsymbol{\Theta}, \mathbf{D})$ has the same form as in Theorem 1:

$$p(\mathbf{y}_n|\boldsymbol{\alpha}_n, \boldsymbol{\Theta}, \mathbf{D}) = \int p(\mathbf{y}_n|\boldsymbol{\alpha}_n, \mathbf{M}_n, \mathbf{D}) p(\mathbf{M}_n|\boldsymbol{\Theta}) d\mathbf{M}_n$$
$$= \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \mathcal{N}(\mathbf{y}_n|\boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}}),$$

where
$$p(\mathbf{y}_n|\boldsymbol{\alpha}_n, \mathbf{M}_n, \mathbf{D}) = \mathcal{N}(\mathbf{y}_n|\sum_j \mathbf{m}_{nj}\alpha_{nj}, \mathbf{D}).$$

The proof is much more complicated (in terms of algebra) and therefore relegated to the supplemental material of the paper.

C. An Example

306

308

311

313

314

315

316

317

322

324

326

329

330

331

334

335

336

338

We give an example to illustrate the basic idea of this paper. Suppose we have M = 4 endmembers with $K_1 = 1$, $K_2 = 2$, $K_3 = 3$, $K_4 = 1$. Their distributions follow (4) with $\mu_{ik}, \Sigma_{ik}, j = 1, 2, 3, 4, k = 1, ..., K_i$. Let the weights of these components be $\pi_{11} = \pi_{41} = 1$, $\pi_{21} = 0.3$, $\pi_{22} = 0.7$, $\pi_{31} = 0.2$, $\pi_{32} = 0.4$, $\pi_{33} = 0.4$. Then, K has 6 entries from the Cartesian product, $\{1\} \times \{1, 2\} \times \{1, 2, 3\} \times \{1\}$. We list the values for $\pi_{\mathbf{k}}$, $\mu_{n\mathbf{k}}$ in Table I. For example, for $\mathbf{k} = (1, 2, 3, 1), \ \pi_{\mathbf{k}} = \pi_{11}\pi_{22}\pi_{33}\pi_{41} = 0.28.$ The value of $\mu_{n\mathbf{k}}$ is a linear combination of μ_{ik} (pick one component for each j) based on the configuration k. Hence, the distribution of y_n in (5) is a Gaussian mixture of 6 components with π_k , $\mu_{n\mathbf{k}}$ given in Table I ($\Sigma_{n\mathbf{k}}$ can be derived similar to $\mu_{n\mathbf{k}}$). Recalling the intuition in Fig. 2, we will show that applying it to hyperspectral unmixing will force each pixel to match all the $\mu_{n\mathbf{k}}$ s, but with emphasis determined by $\pi_{n\mathbf{k}}$.

III. GAUSSIAN MIXTURE MODEL FOR ENDMEMBER VARIABILITY

A. The GMM for Hyperspectral Unmixing

Based on the analysis in Section II, we can model the conditional distribution of all the pixels $\mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_N]^T \in \mathbb{R}^{N \times B}$ given all the abundances $\mathbf{A} := [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_N]^T \in \mathbb{R}^{N \times M}$ ($\boldsymbol{\alpha}_n := [\alpha_{n1}, \dots, \alpha_{nM}]^T$) and GMM parameters, which leads to a maximum *a posteriori* (MAP) problem. Using the result in (5) and assuming the conditional distributions of \mathbf{y}_n are independent, the distribution of \mathbf{Y} given $\mathbf{A}, \boldsymbol{\Theta}, \mathbf{D}$ becomes

$$p(\mathbf{Y}|\mathbf{A}, \mathbf{\Theta}, \mathbf{D}) = \prod_{n=1}^{N} p(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n}, \mathbf{\Theta}, \mathbf{D}).$$
 (15)

Based on the hyperspectral unmixing context, we can set the priors for A. Suppose we use the same prior on A as in [37], i.e.

$$p(\mathbf{A}) \propto \exp\left\{-\frac{\beta_1}{2} \operatorname{Tr}\left(\mathbf{A}^T \mathbf{L} \mathbf{A}\right) + \frac{\beta_2}{2} \operatorname{Tr}\left(\mathbf{A}^T \mathbf{A}\right)\right\}$$
$$= \exp\left\{-\frac{\beta_1}{2} \operatorname{Tr}\left(\mathbf{A}^T \mathbf{K} \mathbf{A}\right)\right\}, \tag{16}$$

where **L** is a *graph Laplacian* matrix constructed from w_{nm} , n, m = 1, ..., N with $w_{nm} = e^{-\|\mathbf{y}_n - \mathbf{y}_m\|^2/2B\eta^2}$ for neighboring pixels and 0 otherwise. We have $\operatorname{Tr}(\mathbf{A}^T \mathbf{L} \mathbf{A}) = \frac{1}{2} \sum_{n,m} w_{nm} \|\boldsymbol{\alpha}_n - \boldsymbol{\alpha}_m\|^2$, $\mathbf{K} = \mathbf{L} - \frac{\beta_2}{\beta_1} \mathbf{I}_N$ (suppose $\beta_1 \neq 0$) with β_1 controlling smoothness and β_2 controlling sparsity of the abundance maps.

From the conditional density function and the priors, Bayes' theorem says the posterior is given by

$$p(\mathbf{A}, \mathbf{\Theta}|\mathbf{Y}, \mathbf{D}) \propto p(\mathbf{Y}|\mathbf{A}, \mathbf{\Theta}, \mathbf{D}) p(\mathbf{A}) p(\mathbf{\Theta}),$$
 (17)

where $p(\mathbf{\Theta})$ is assumed to follow a uniform distribution. Maximizing $p(\mathbf{A}, \mathbf{\Theta}|\mathbf{Y}, \mathbf{D})$ is equivalent to minimizing $-\log p(\mathbf{A}, \mathbf{\Theta}|\mathbf{Y}, \mathbf{D})$, which reduces to the following form by combining (5), (15), (16) and (17):

$$\mathcal{E}(\mathbf{A}, \mathbf{\Theta}) = -\sum_{n=1}^{N} \log \sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \mathcal{N}(\mathbf{y}_{n} | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}}) + \mathcal{E}_{\text{prior}}(\mathbf{A}),$$

s.t.
$$\pi_{\mathbf{k}} \ge 0$$
, $\sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} = 1$, $\alpha_{nj} \ge 0$, $\sum_{j=1}^{M} \alpha_{nj} = 1$, $\forall n$ (18)

where $\mathcal{E}_{prior}(\mathbf{A}) = \frac{\beta_1}{2} \text{Tr} (\mathbf{A}^T \mathbf{K} \mathbf{A})$, and $\boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}}$ are defined in (6).

B. Relationships to Least-Squares, NCM, and MESMA

Let us focus on the first term in (18) and call it the *data fidelity term*. We can relate it to NCM and the least-squares term $\sum_n \|\mathbf{y}_n - \sum_j \alpha_{nj} \mathbf{m}_j\|^2$ as used in previous research. The data fidelity term in NCM follows (3) and is based on minimizing the negative log-likelihood

$$-\log p\left(\mathbf{Y}\right) = -\log \prod_{n=1}^{N} p\left(\mathbf{y}_{n}\right) = -\sum_{n=1}^{N} \log \mathcal{N}\left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n1}, \boldsymbol{\Sigma}_{n1}\right)$$
(19)

by assuming \mathbf{y}_n s are independent, where $\boldsymbol{\mu}_{n1} := \sum_j \alpha_{nj} \boldsymbol{\mu}_j$, $\boldsymbol{\Sigma}_{n1} := \sum_j \alpha_{nj}^2 \boldsymbol{\Sigma}_j + \sigma^2 \mathbf{I}_B$. Expanding (19) using the form of the Gaussian distribution leads to the objective function

$$\sum_{n=1}^{N} \log |\mathbf{\Sigma}_{n1}| + \sum_{n=1}^{N} (\mathbf{y}_{n} - \boldsymbol{\mu}_{n1})^{T} \mathbf{\Sigma}_{n1}^{-1} (\mathbf{y}_{n} - \boldsymbol{\mu}_{n1}). \quad (20)$$

We can see that the least-squares minimization is a special case of NCM with $\|\mathbf{\Sigma}_j\|_F \to 0$, i.e. when there is little endmember variability.

The proposed GMM further generalizes NCM from a statistical perspective. Since π_{jk} represents the prior probability of the latent variable in a GMM, $\pi_{\mathbf{k}}$ represents the prior probability of picking a combination. If we see \mathbf{k} as a (discrete) random variable whose sample space is \mathcal{K} , (5) can be seen as

$$p\left(\mathbf{y}_{n}|\mathbf{\alpha}_{n},\mathbf{\Theta},\mathbf{D}\right) = \sum_{\mathbf{k}\in\mathcal{K}} p\left(\mathbf{k}\right) p\left(\mathbf{y}_{n}|\mathbf{k},\mathbf{\alpha}_{n},\mathbf{\Theta},\mathbf{D}\right),$$
 38

where $p(\mathbf{k}) = \pi_{\mathbf{k}}$ and $p(\mathbf{y}_n | \mathbf{k}, \alpha_n, \boldsymbol{\Theta}, \mathbf{D}) = \mathcal{N}(\mathbf{y}_n | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}})$. From this perspective, each pixel is generated by first sampling \mathbf{k} , then sampling a Gaussian distribution determined by \mathbf{k} , $\boldsymbol{\Theta}$. Unlike NCM that tries to make each \mathbf{y}_n close to $\boldsymbol{\mu}_{n\mathbf{1}}$ which is a linear combination of a fixed set $\{\boldsymbol{\mu}_j\}$, GMM further generalizes it by trying to make \mathbf{y}_n close to every $\boldsymbol{\mu}_{n\mathbf{k}}$ which are all the possible linear combinations of $\{\boldsymbol{\mu}_{jk}\}$. It makes sense that the summation in (18) is weighted by $\pi_{\mathbf{k}}$ in a way that if one combination has a high probability to appear, i.e. $\pi_{\mathbf{k}}$ is larger for a certain \mathbf{k} , the effort is biased to make \mathbf{y}_n closer to this particular $\boldsymbol{\mu}_{n\mathbf{k}}$. Fig. 2 shows the differences among these.

The widely adopted MESMA takes a library of endmember spectra as input, tries all the combinations and pick the combination with least reconstruction error. The philosophy is similar to our model despite the fundamental difference that MESMA is explicit whereas we are implicit in terms of linear combinations. Compared to MESMA, the GMM approach separates the library into M groups where each group represents a material and is clustered into several centers, such that the combination can only take place by picking one center from each group. Also, the size of each cluster affects the probability of picking its center. Hence, our model can adapt to very large library sizes as long as the number of clusters does not increase too much.

C. Optimization

Estimating the parameters of GMMs has been studied extensively, from early expectation maximization (EM) from the statistical community to projection based clustering from the computer science community [38], [39]. There are simple and deterministic algorithms, which usually require the centers of Gaussian be separable. However, we face a more challenging problem since each pixel is generated by a different GMM determined by the coefficients α_n . Since EM can be seen as a special case of Majoriziation-Minimization algorithms [40], which is more flexible, we adopt this approach. Considering that we have too many parameters \mathbf{A} , $\mathbf{\Theta}$ to update in the M step, they are updated sequentially as long as the complete data log-likelihood increases. This is also called *generalized* expectation maximization (GEM) [41].

Following the routine of EM, the E step calculates the posterior probability of the latent variable given the observed data and old parameters

$$\gamma_{n\mathbf{k}} = \frac{\pi_{\mathbf{k}} \mathcal{N}\left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}}\right)}{\sum_{\mathbf{k} \in \mathcal{K}} \pi_{\mathbf{k}} \mathcal{N}\left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}}\right)}.$$
 (21)

The M step usually maximizes the expected value of the complete data log-likelihood. Here, we have priors in the Bayesian formulation. Hence, we need to minimize

$$\mathcal{E}_{M} = -\sum_{n=1}^{N} \sum_{\mathbf{k} \in \mathcal{K}} \gamma_{n\mathbf{k}} \left\{ \log \pi_{\mathbf{k}} + \log \mathcal{N} \left(\mathbf{y}_{n} | \boldsymbol{\mu}_{n\mathbf{k}}, \boldsymbol{\Sigma}_{n\mathbf{k}} \right) \right\} + \mathcal{E}_{\text{prior}}.$$

This leads to a common update step for $\pi_{\mathbf{k}}$ as

$$\pi_{\mathbf{k}} = \frac{1}{N} \sum_{n=1}^{N} \gamma_{n\mathbf{k}}.$$
 (23)

We now focus on updating $\{\mu_{jk}, \Sigma_{jk}\}$ and **A**. To achieve this, we require the derivatives of \mathcal{E}_M in (22) w.r.t. $\mu_{jk}, \Sigma_{jk}, \alpha_{nj}$. After some tedious algebra using (6), we get

$$\frac{\partial \mathcal{E}_{M}}{\partial \boldsymbol{\mu}_{jl}} = -\sum_{n=1}^{N} \sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_{j}} \alpha_{nj} \boldsymbol{\lambda}_{n\mathbf{k}}$$
(24) 437

$$\frac{\partial \mathcal{E}_M}{\partial \mathbf{\Sigma}_{jl}} = -\sum_{n=1}^N \sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \alpha_{nj}^2 \mathbf{\Psi}_{n\mathbf{k}},\tag{25}$$

$$\frac{\partial \mathcal{E}_{M}}{\partial \alpha_{nj}} = -\sum_{\mathbf{k} \in \mathcal{K}} \boldsymbol{\lambda}_{n\mathbf{k}}^{T} \boldsymbol{\mu}_{jk_{j}} - 2\alpha_{nj} \sum_{\mathbf{k} \in \mathcal{K}} \operatorname{Tr} \left(\boldsymbol{\Psi}_{n\mathbf{k}}^{T} \boldsymbol{\Sigma}_{jk_{j}} \right) + \beta_{1} \left(\mathbf{K} \mathbf{A} \right)_{nj}, \tag{26}$$

where $\lambda_{n\mathbf{k}} \in \mathbb{R}^{B \times 1}$ and $\Psi_{n\mathbf{k}} \in \mathbb{R}^{B \times B}$ are given by

$$\lambda_{n\mathbf{k}} = \gamma_{n\mathbf{k}} \mathbf{\Sigma}_{n\mathbf{k}}^{-1} \left(\mathbf{y}_n - \boldsymbol{\mu}_{n\mathbf{k}} \right), \tag{27}$$

$$\Psi_{n\mathbf{k}} = \frac{1}{2} \gamma_{n\mathbf{k}} \mathbf{\Sigma}_{n\mathbf{k}}^{-T} \left(\mathbf{y}_n - \boldsymbol{\mu}_{n\mathbf{k}} \right) \left(\mathbf{y}_n - \boldsymbol{\mu}_{n\mathbf{k}} \right)^T \mathbf{\Sigma}_{n\mathbf{k}}^{-T} - \frac{1}{2} \gamma_{n\mathbf{k}} \mathbf{\Sigma}_{n\mathbf{k}}^{-T}.$$
(28)

It is better to represent the derivatives in matrix forms for the sake of implementation convenience. Considering the multiple summations in (24), (25) and (26), we can write them as

$$\frac{\partial \mathcal{E}_M}{\partial \boldsymbol{\mu}_{jl}} = -\sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \left(\mathbf{A}^T \boldsymbol{\Lambda}_{\mathbf{k}} \right)_j, \tag{29}$$

$$\frac{\partial \mathcal{E}_M}{\partial \text{vec}\left(\mathbf{\Sigma}_{jl}\right)} = -\sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \left((\mathbf{A} \circ \mathbf{A})^T \, \mathbf{\Psi}_{\mathbf{k}} \right)_j, \tag{30}$$

$$\frac{\partial \mathcal{E}_M}{\partial \mathbf{A}} = -\sum_{\mathbf{k} \in \mathcal{K}} \mathbf{\Lambda}_{\mathbf{k}} \mathbf{R}_{\mathbf{k}}^T - 2\mathbf{A} \circ \sum_{\mathbf{k} \in \mathcal{K}} \mathbf{\Psi}_{\mathbf{k}} \mathbf{S}_{\mathbf{k}}^T + \beta_1 \mathbf{K} \mathbf{A}, \quad (31)$$

where $\Lambda_{\mathbf{k}} \in \mathbb{R}^{N \times B}$, $\Psi_{\mathbf{k}} \in \mathbb{R}^{N \times B^2}$ denote the matrices formed by $\{\lambda_{n\mathbf{k}}, \Psi_{n\mathbf{k}}\}$ as follows

$$\Lambda_{\mathbf{k}} := [\lambda_{1\mathbf{k}}, \lambda_{2\mathbf{k}}, \dots, \lambda_{N\mathbf{k}}]^{T},$$

$$\Psi_{\mathbf{k}} := [\text{vec}(\Psi_{1\mathbf{k}}), \text{vec}(\Psi_{2\mathbf{k}}), \dots, \text{vec}(\Psi_{N\mathbf{k}})]^{T},$$
453

and $\mathbf{R}_{\mathbf{k}} \in \mathbb{R}^{M \times B}$, $\mathbf{S}_{\mathbf{k}} \in \mathbb{R}^{M \times B^2}$ are defined by

$$\mathbf{R}_{\mathbf{k}} := \left[\mu_{1k_1}, \mu_{2k_2}, \dots, \mu_{Mk_M} \right]^T, \tag{32}$$

$$\mathbf{S}_{\mathbf{k}} := \left[\operatorname{vec} \left(\mathbf{\Sigma}_{1k_1} \right), \operatorname{vec} \left(\mathbf{\Sigma}_{2k_2} \right), \dots, \operatorname{vec} \left(\mathbf{\Sigma}_{Mk_M} \right) \right]^T.$$
 (33)

The minimum of \mathcal{E}_M corresponds to $\frac{\partial \mathcal{E}_M}{\partial \mu_{jl}} = 0$, $\frac{\partial \mathcal{E}_M}{\partial \Sigma_{jl}} = 0$, and $\frac{\partial \mathcal{E}_M}{\partial \mathbf{A}} = 0$ if the optimization problem is unconstrained. However, since we have the non-negativity and sum-to-one constraint to α_{nj} and positive definite constraint of Σ_{jk} , minimizing \mathcal{E}_M is very difficult. Therefore, in each M step, we only decrease this objective function by *projected gradient descent* (please see [42 and 43, Sec. 2.3]) using (29), (30) and (31), where the projection functions for \mathbf{A} and $\{\Sigma_{jk}\}$ are the same as in [37].

Finally, from the estimated $\pi_{\mathbf{k}}$, we can recover the sets of weights as $\pi_{jl} = \sum_{\mathbf{k} \in \mathcal{K}} \delta_{lk_j} \pi_{\mathbf{k}}$.

D. Model Selection

(22)

The number of components K_j can be specified or estimated from the data. For the latter case, we have some pure pixels and estimate K_j by deploying a standard

model selection method. Suppose we have N_j pure pixels $\mathbf{Y}_j := \left[\mathbf{y}_1^j, \mathbf{y}_2^j, \dots, \mathbf{y}_{N_j}^j\right]^T \in \mathbb{R}^{N_j \times B}$ for the jth endmember, $f_{\mathbf{m}_j}\left(\mathbf{y}|\mathbf{\Theta}_j\right)$ is the estimated density function with $\mathbf{\Theta}_j := \left\{\pi_{jk}, \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk} : k = 1, \dots, K_j\right\}$, $g_{\mathbf{m}_j}\left(\mathbf{y}\right)$ is the true density function. The information criterion based model selection approach tries to find K_j that minimizes their difference, e.g. the Kullback-Leibler (KL) divergence

$$\begin{split} \mathcal{D}_{\mathrm{KL}}\left(g_{\mathbf{m}_{j}} \| f_{\mathbf{m}_{j}}\right) &= \int_{\mathbb{R}^{B}} g_{\mathbf{m}_{j}}\left(\mathbf{y}\right) \log \frac{g_{\mathbf{m}_{j}}\left(\mathbf{y}\right)}{f_{\mathbf{m}_{j}}\left(\mathbf{y} | \mathbf{\Theta}_{j}\right)} d\mathbf{y} \\ &\approx -\frac{1}{N_{j}} \sum_{i}^{N_{j}} \log f_{\mathbf{m}_{j}}\left(\mathbf{y}_{n}^{j} | \mathbf{\Theta}_{j}\right) + \mathrm{const}, \end{split}$$

where the approximation of $\int g_{\mathbf{m}_j}(\mathbf{y}) \log f_{\mathbf{m}_j}(\mathbf{y}|\mathbf{\Theta}_j) d\mathbf{y}$ by the log-likelihood is usually biased as the empirical distribution function is closer to the fitted distribution than the true one. Akaike's information criterion is one way to approximate the bias. Here, we use the cross-validation-based information criterion (CVIC) to correct for the bias [44], [45]. Let

$$\mathcal{L}_{\mathbf{Y}_{j}}\left(\mathbf{\Theta}_{j}\right) = \sum_{n=1}^{N_{j}} \log f_{\mathbf{m}_{j}}\left(\mathbf{y}_{n}^{j} | \mathbf{\Theta}_{j}\right). \tag{34}$$

The *V*-fold cross validation (we use V = 5 here) divides the input set \mathbf{Y}_j into V subsets $\left\{\mathbf{Y}_j^1, \mathbf{Y}_j^2, \dots, \mathbf{Y}_j^V\right\}$ with equal sizes. Then for each subset \mathbf{Y}_j^v , $v = 1, \dots, V$, the remaining data are used to replace \mathbf{Y}_j in (34) such that (34) is maximized by $\mathbf{\Theta}_j^v$. Then $\mathcal{L}_{K_j} = \sum_v \mathcal{L}_{\mathbf{Y}_j^v} \left(\mathbf{\Theta}_j^v\right)$ is evaluated and the optimal $\hat{K}_j = \arg\max_{K_j} \mathcal{L}_{K_j}$.

E. Implementation Details

The algorithm can be implemented in a supervised or unsupervised manner. In both cases, because of the large computational cost, we project the pixel data to a low dimensional space by principal component analysis (PCA) and perform the optimization, the result then projected back to the original space. Let $\mathbf{E} \in \mathbb{R}^{B \times d}$ be the projection matrix and $\mathbf{c} \in \mathbb{R}^{B}$ be the translation vector, then

$$\mathbf{E}^{T}\left(\mathbf{y}_{n}-\mathbf{c}\right)=\sum_{j=1}^{M}\mathbf{E}^{T}\left(\mathbf{m}_{nj}-\mathbf{c}\right)\alpha_{nj}+\mathbf{E}^{T}\mathbf{n}_{n}.$$

This means that for the projected pixels, the *j*th endmember $\mathbf{m}'_{nj} = \mathbf{E}^T (\mathbf{m}_{nj} - \mathbf{c})$ follows a distribution

$$p\left(\mathbf{m}_{nj}'|\mathbf{\Theta}\right) = \sum_{k=1}^{K_j} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj}'|\mathbf{E}^T\left(\boldsymbol{\mu}_{jk} - \mathbf{c}\right), \mathbf{E}^T \boldsymbol{\Sigma}_{jk} \mathbf{E}\right)$$

and the noise $\mathbf{n}'_n = \mathbf{E}^T \mathbf{n}_n$ follows $\mathcal{N}(\mathbf{n}'_n | \mathbf{0}, \mathbf{E}^T \mathbf{D} \mathbf{E})$.

In the supervised unmixing scenario, we assume that a library of endmember spectra is known. After estimating the number of components following Section III-D, and calculating Θ using the standard EM algorithm, we only need to update $\gamma_{n\mathbf{k}}$ by (21) and \mathbf{A} by (31) with $\pi_{\mathbf{k}}$, μ_{jk} and Σ_{jk} fixed. The initialization of \mathbf{A} can utilize the multiple combinations of means. For each α_n , we first set

 $\alpha_{n\mathbf{k}} \leftarrow (\mathbf{R}_{\mathbf{k}}\mathbf{R}_{\mathbf{k}}^T + \epsilon \mathbf{I}_M)^{-1} \mathbf{R}_{\mathbf{k}}\mathbf{y}_n$, then project it to the simplex space, and finally set $\alpha_n \leftarrow \alpha_{n\mathbf{Q}}$ with $\hat{\mathbf{k}} = \arg\min_{\mathbf{k}} \|\mathbf{y}_n - \mathbf{R}_{\mathbf{k}}^T \alpha_{n\mathbf{k}}\|^2$, i.e. choose the $\alpha_{n\mathbf{k}}$ that minimizes the reconstruction error.

In the unsupervised unmixing scenario, we will assume the resolution is high enough such that the hyperspectral image can be segmented into several regions where the interior pixels in each region are pure pixels. The optimization is performed in several steps, where we first obtain a segmentation result, then use CVIC to determine the number of components, and finally estimate $\bf A$ with $\bf \Theta$ fixed. The details are given as follows.

Step 1: Initialization. We start with $K_j = 1$, $\forall j$ and use K-means to find the initial means $\mathbf{R_1}$. The initial \mathbf{A} is set to $\mathbf{A} \leftarrow \mathbf{Y}\mathbf{R_1}^T \left(\mathbf{R_1}\mathbf{R_1}^T + \epsilon \mathbf{I}_M\right)^{-1}$ (by minimizing $\|\mathbf{Y} - \mathbf{A}\mathbf{R_1}\|_F^2$), then projected to the valid simplex space as in [37]. The initial covariance matrices are set to $\Sigma_{j1} \leftarrow 0.1^2 \mathbf{I}_B$, $\forall j$. For the noise matrix \mathbf{D} , although there is research focused on noise estimation [46], [47], endmember variability was not considered and validation was performed only for the simple LMM assumption. Hence, we use an empirical value $\mathbf{D} = 0.001^2 \mathbf{I}_B$, which is usually much less than the variability of covariance matrices in (6).

Step 2: Segmentation. Given the initial conditions, we use the GEM algorithm to iteratively update $\gamma_{n\mathbf{k}}$ by (21), $\pi_{\mathbf{k}}$ by (23), μ_{jk} by (29), \mathbf{A} by (31) while keeping Σ_{jk} fixed. For $\gamma_{n\mathbf{k}}$ and $\pi_{\mathbf{k}}$, a direct update equation is available. For μ_{jk} , we can use gradient descent. For \mathbf{A} , since we have the non-negativity and sum-to-one constraints, a projected gradient descent similar to the one used in [37] can be applied. To ensure a segmentation effect, a large β_2 is used in this step.

Step 3: Model selection and abundance estimation. Using the segmentation-like abundance maps from the previous step, we can obtain the interior pixels \mathbf{Y}_j (assumed pure) by thresholding the abundances (e.g. $\alpha_{nj} > 0.99$) and performing image erosion to trim the boundaries with structure element size r_{se} (can be decreased gradually if large enough to trim all the pixels). Following Section III-D, we can determine the number of components K_j and further calculate $\mathbf{\Theta}_j$ by standard EM. Since β_2 is relatively large in the previous step, it is reduced by $\beta_2 \leftarrow \zeta \beta_2$ where $\zeta = 0.05$. Then we restart the optimization to estimate the abundances with $\mathbf{\Theta}$ fixed.

F. Complexity Analysis

The abundance estimation algorithm is an iterative process. Since we used projected gradient descent with adaptive step sizes, the number of iterations is usually not large as shown in [43] and [48]. For each iteration, it starts with calculating $\mu_{n\mathbf{k}}$ and $\Sigma_{n\mathbf{k}}$ in (6), where storing all $\mu_{n\mathbf{k}}$ ($\Sigma_{n\mathbf{k}}$) requires $O(|\mathcal{K}|NB)$ ($O(|\mathcal{K}|NB^2)$), the computation takes $O(|\mathcal{K}|NMB)$ ($O(|\mathcal{K}|NMB^2)$). Suppose the Cholesky factorization and the matrix inversion of a B by B matrix both take $O(B^3)$ time, and $N\gg B>M$. Evaluating $\log \mathcal{N}(\mathbf{y}_n|\mu_{n\mathbf{k}},\Sigma_{n\mathbf{k}})$ by the Cholesky factorization will take $O(B^3)$, hence updating all the $\gamma_{n\mathbf{k}}$ takes $O(|\mathcal{K}|NB^3)$, which is also the required time for evaluating the objective function (18). The calculation of $\lambda_{n\mathbf{k}}$, $\Psi_{n\mathbf{k}}$ (in (27) and (28)) will be

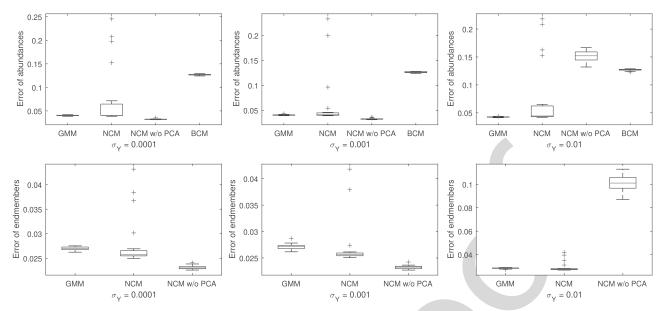


Fig. 3. Abundance and endmember error statistics from 20 synthetic images for each noise level in the supervised unmixing scenario.

dominated by the inversion of $\Sigma_{n\mathbf{k}}$ which takes $O(B^3)$, hence the overall calculation takes $O(|\mathcal{K}|NB^3)$ with storage the same as $\mu_{n\mathbf{k}}$ and $\Sigma_{n\mathbf{k}}$. Then if we move to calculating the derivatives in (29), (30) and (31), it is easy to verify that the computational costs are $O(|\mathcal{K}|NMB)$, $O(|\mathcal{K}|NMB^2)$, $O(|\mathcal{K}|NMB^2)$ respectively (Note that \mathbf{K} is a banded matrix so the computation involving it is linear). Reviewing the above process, we conclude that the spatial complexity is dominated by $O(|\mathcal{K}|NB^2)$ and the time complexity is dominated by $O(|\mathcal{K}|NB^3)$.

G. Estimation of Endmembers for Each Pixel

While the previous sections discuss the estimation of the abundances and endmember distribution parameters, they do not actually estimate the endmembers $\{\mathbf{m}_{nj}: n=1,\ldots,N, j=1,\ldots,M\}$ for each pixel. In this Section, we will discuss this additional problem and note its absence in the previous NCM literature.

Theorem 2 implies that we can view the proposed conditional density (5) as modeling the noise as a Gaussian random variable followed by marginalizing over \mathbf{M}_n , which is usually achieved by the evidence approximation in the machine learning literature due to the intractability of the integral ([49, Sec. 3.5]). Since we have \mathbf{A} , $\mathbf{\Theta}$ obtained from the previous Sections, we can get the posterior of \mathbf{M}_n from this model:

$$p\left(\mathbf{M}_{n}|\mathbf{y}_{n},\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right) \propto p\left(\mathbf{y}_{n},\mathbf{M}_{n}|\boldsymbol{\alpha}_{n},\boldsymbol{\Theta},\mathbf{D}\right)$$
$$= p\left(\mathbf{y}_{n}|\boldsymbol{\alpha}_{n},\mathbf{M}_{n},\mathbf{D}\right) p\left(\mathbf{M}_{n}|\boldsymbol{\Theta}\right). \quad (35)$$

Maximizing $\log p(\mathbf{M}_n|\mathbf{y}_n, \boldsymbol{\alpha}_n, \boldsymbol{\Theta}, \mathbf{D})$ gives us another minimization problem

$$\mathcal{E}\left(\mathbf{M}_{n}\right) = \frac{1}{2} \left(\mathbf{y}_{n} - \mathbf{M}_{n}^{T} \boldsymbol{\alpha}_{n}\right)^{T} \mathbf{D}^{-1} \left(\mathbf{y}_{n} - \mathbf{M}_{n}^{T} \boldsymbol{\alpha}_{n}\right)$$
$$- \sum_{i=1}^{M} \log \sum_{k=1}^{K_{j}} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right) \quad (36)$$

obtained by plugging (13) and (14) into (35). Note that this objective function has an intuitive interpretation as the first term minimizes the reconstruction error while the second term forces the endmembers close to the centers of each GMM. The weight factor between the two terms is the noise. From an algebraic perspective, since there are also logarithms of sums of Gaussian functions in this objective, we can also use the EM algorithm for ease of optimization. In the E step, the soft membership is calculated by

$$\gamma_{njk} = \frac{\pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right)}{\sum_{k} \pi_{jk} \mathcal{N}\left(\mathbf{m}_{nj} | \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\right)}, \quad k = 1, \dots, K_{j}.$$

In the M step, the derivative w.r.t. \mathbf{m}_{nj} is obtained as

$$\frac{\partial \mathcal{E}}{\partial \mathbf{m}_{nj}} = -\mathbf{D}^{-1} \left(\mathbf{y}_n - \mathbf{M}_n^T \boldsymbol{\alpha}_n \right) \alpha_{nj}$$

$$\mathbf{\Sigma}^{-1} \left(\mathbf{y}_n - \mathbf{M}_n^T \boldsymbol{\alpha}_n \right)$$
613

$$+\sum_{k=1}^{K_j}\gamma_{njk}\boldsymbol{\Sigma}_{jk}^{-1}\left(\mathbf{m}_{nj}-\boldsymbol{\mu}_{jk}\right).$$

Instead of deploying gradient descent in the M step for estimating the abundances, combining the derivatives for all j actually leads to a closed form solution

$$\operatorname{vec}\left(\mathbf{M}_{n}^{T}\right) = \left\{\boldsymbol{\alpha}_{n}\boldsymbol{\alpha}_{n}^{T} \otimes \mathbf{D}^{-1} + \operatorname{diag}\left(\mathbf{C}_{n1}, \dots, \mathbf{C}_{nM}\right)\right\}^{-1}$$

$$\left\{\operatorname{vec}\left(\mathbf{D}^{-1}\mathbf{y}_{n}\boldsymbol{\alpha}_{n}^{T}\right) + \mathbf{d}_{n}\right\}$$
618

where $\mathbf{C}_{nj} \in \mathbb{R}^{B \times B}$ and $\mathbf{d}_n := (\mathbf{d}_{n1}^T, \dots, \mathbf{d}_{nM}^T)^T \in \mathbb{R}^{MB \times 1}$ are defined as

$$\mathbf{C}_{nj} := \sum_{k=1}^{K_j} \gamma_{njk} \mathbf{\Sigma}_{jk}^{-1}, \ \mathbf{d}_{nj} := \sum_{k=1}^{K_j} \gamma_{njk} \mathbf{\Sigma}_{jk}^{-1} \boldsymbol{\mu}_{jk}.$$

In practice, despite the need to estimate a large $M \times B \times N$ tensor, the time cost is actually much less than the estimation of abundances because of the closed form update equation in the M step. An interesting fact is that γ_{njk} measures the closeness of estimated endmembers to clusters centers, hence

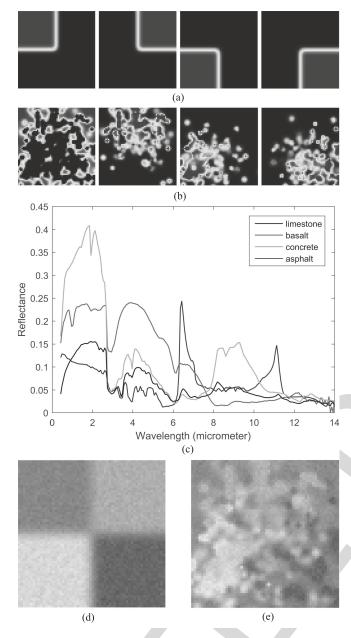


Fig. 4. Unsupervised synthetic dataset. (a) and (b) are abundance maps for two images. (c) shows original spectra from the ASTER library. (d) and (e) show the color images.

may provide a clue on which cluster is sampled to generate an endmember.

IV. RESULTS

In the following experiments, we implemented the algorithm in MATLAB® and compared the proposed GMM with NCM, BCM (spectral version with quadratic programming) [15] on synthetic and real images. As mentioned previously, for GMM, the original image data were projected to a subspace with 10 dimensions to speed up the computation for abundance estimation. NCM was implemented as a supervised algorithm wherein we input the ground truth pure pixels (in the image

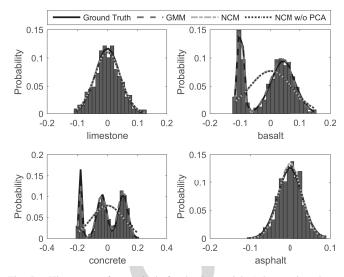


Fig. 5. Histograms of pure pixels for the 4 materials (when projected to a 1-dimensional space determined by performing PCA on the pure pixels of each material) and the ground truth and estimated distributions (also projected to the same direction) for the first image of the unsupervised synthetic dataset. The probability of each distribution is calculated by multiplying the value of the density function at each bin location with the bin size.

TABLE II $L_2 \ \, \text{DISTANCE Between the Fitted Distributions (GMM, NCM)} \\ \text{AND THE GROUND TRUTH DISTRIBUTIONS FOR THE FIRST} \\ \text{IMAGE OF THE UNSUPERVISED SYNTHETIC DATASET}$

	$\times 10^{6}$	Limestone	Basalt	Concrete	Asphalt	Mean
1	GMM	4.45	3.46	3.41	4.28	3.85
Ì	NCM	4.27	5.86	4.95	4.02	4.77

TABLE III
ABUNDANCE ERRORS FOR THE UNSUPERVISED SYNTHETIC DATASET

	$\times 10^{-4}$	GMM	NCM	NCM w/o PCA	BCM
	Limestone	50	107	92	126
	Basalt	40	74	67	158
Image	Concrete	41	66	62	186
l m	Asphalt	69	141	123	292
	Mean	59	97	86	190
	Limestone	157	1086	396	231
e 2	Basalt	126	445	270	204
Image	Concrete	103	985	229	206
Im	Asphalt	225	170	706	445
	Mean	153	671	400	272

with extreme abundances), modeled them by Gaussian distributions, and obtained the abundance maps by maximizing the log-likelihood. We considered two versions of NCM, one in the same subspace as GMM (referred to as NCM), the other in the original spectral space (referred to as NCM without PCA). Since BCM is also a supervised unmixing algorithm, ground truth pure pixels were again taken as input and the results were the abundance maps. For GMM and the two versions of NCM, using the algorithm in Section III-G we can obtain the endmembers for each pixel. All the parameters of GMM (except the structure element size r_{se}) were set to $\beta_1 = 5$, $\beta_2 = 5$ unless specified throughout the experiments.

For comparison of endmember distributions, we calculated the L_2 distance $\left(\int |f(\mathbf{x}) - g(\mathbf{x})|^2 d\mathbf{x}\right)^{1/2}$ between the fitted distribution and the ground truth one, where the latter was only

¹The code of GMM is available on GitHub (https://github.com/zhouyuanzxcv/Hyperspectral).

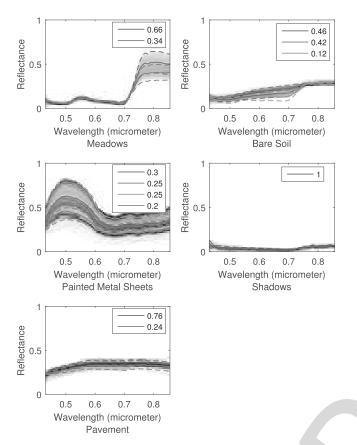


Fig. 6. Estimated GMM in the wavelength-reflectance space for the Pavia University dataset. The background gray image represents the histogram created by placing the pure pixel spectra into the reflectance bins at each wavelength. The different colors represent different components, where the solid curve is the center μ_{jk} , the dashed curves are $\mu_{jk} \pm 2\sigma_{jk} \mathbf{v}_{jk}$ (σ_{jk} is the square root of the large eigenvalue of Σ_{jk} while \mathbf{v}_{jk} is the corresponding eigenvector), and the legend shows the prior probabilities.

available for the synthetic dataset. For comparison of abundances, we calculated the root mean squared error (RMSE) $\left(\frac{1}{N}\sum_{n}|\alpha_{nj}^{GT}-\alpha_{nj}^{est}|^2\right)^{1/2}$ where α_{nj}^{GT} are the ground truth abundances and α_{nj}^{est} are the estimated values. Since only some pure pixels were identified as ground truth in the real datasets, we calculated error $j=\left(\frac{1}{|\mathcal{I}|}\sum_{n\in\mathcal{I}}|\alpha_{nj}^{GT}-\alpha_{nj}^{est}|^2\right)^{1/2}$ given the pure pixel index set \mathcal{I} . For comparison of endmembers, the same error formula and overall schema were used, i.e. for an index set \mathcal{I}_j of pure pixels for the jth endmember (in the real datasets), error $j=\frac{1}{|\mathcal{I}_j|}\sum_{n\in\mathcal{I}_j}\left(\frac{1}{B}\|\mathbf{m}_{nj}^{GT}-\mathbf{m}_{nj}^{est}\|^2\right)^{1/2}$.

A. Synthetic Datasets

The algorithms were tested for two cases of synthetic images, a supervised case and an unsupervised case.

1) Supervised: In this case, a library of ground truth endmembers were input and the abundances were estimated. The images were of size 60×60 with 103 wavelengths from 430 nm to 860 nm (≤ 5 nm spectral resolution) and created with two endmember classes, meadows and painted metal sheets, whose spectra were drawn randomly from the ground truth of the Pavia University dataset (shown in Fig. 1, meadows have 309 samples and painted metal sheets have

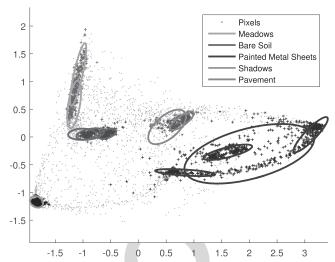


Fig. 7. Scatter plot of the Pavia University dataset with the estimated GMM. The gray dots are the projected pixels by PCA. The darkened dots with a color represent the ground truth pure pixels for a material. The ellipses with the same color represent the projected Gaussian components (twice the standard deviation along the major and minor axes, covering 86% of the total probability mass) for one endmember.

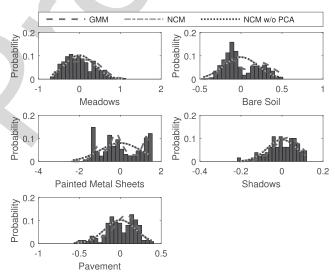


Fig. 8. Histograms of pure pixels for the Pavia University dataset and the estimated distributions from GMM and NCM when projected to 1 dimension.

941 samples in the ROI). Since painted metal sheets have multiple modes in the distribution, it should reflect a true difference between GMM and the other distributions. The abundances were sampled from a Dirichlet distribution so each pixel had random values. Also, an additive noise sampled from $\mathcal{N}(\mathbf{n}_n|\mathbf{0},\mathbf{D})$ was added to the mixed spectra, where the noise was assumed to be independent at different wavelengths, i.e. $\mathbf{D} = \operatorname{diag}\left(\sigma_1^2,\ldots,\sigma_B^2\right)$ while σ_k was again sampled from a uniform distribution on $[0,\sigma_V]$.

We tested the algorithms for different σ_Y . The effects of priors were all removed in this case, i.e. $\beta_1 = 0$, $\beta_2 = 0$. Fig. 3 shows the box plots of abundance and endmember errors. We can see that GMM has small errors in general for different noise levels. NCM also has relatively small errors in most cases, but tends to produce large errors occasionally

696

698

700

702

704

706

707

709

711

712

713

715

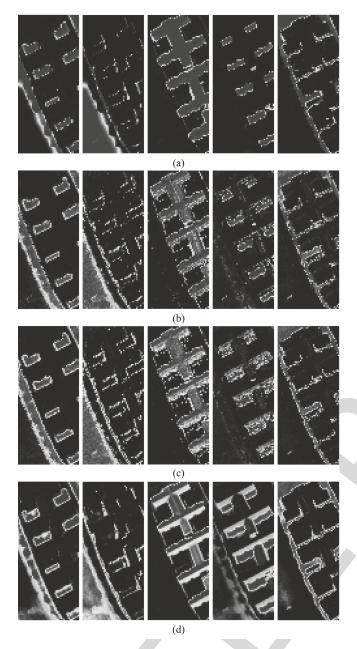


Fig. 9. Abundance maps for the Pavia University dataset. The corresponding endmembers from left to right are meadows, bare soil, painted metal sheets, shadows and pavement. (a) GMM. (b) NCM. (c) NCM w/o PCA. (d) BCM.

TABLE IV
ABUNDANCE AND ENDMEMBER ERRORS FOR PAVIA UNIVERSITY

$\times 10^{-4}$	GMM	NCM	NCM w/o PCA	BCM
Meadow	187 \ 44 ^a	405 \ 113	378 \ 114	711
Soil	175 \ 30	581 \ 68	507 \ 66	1049
Metal	476 \ 49	1236 \ 237	917 \ 349	1285
Shadow	44 \ 44	736 \ 48	914 \ 34	1287
Pavement	473 \ 39	1064 \ 114	333 \ 103	612
Mean	271 \ 41	804 \ 116	610 \ 133	989

^a the numbers in ".\." denote the abundance and endmember errors.

(4 out of 20 runs). NCM without PCA has very good results except for large noise, where it performed worst among all the methods. BCM has the largest errors overall. For the endmembers, although NCM or NCM without PCA sometimes has

691

693

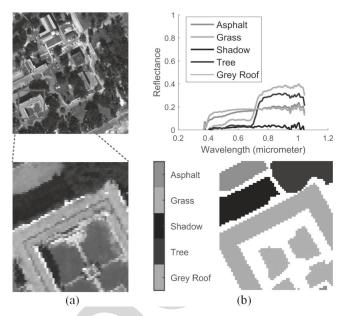


Fig. 10. (a) Original RGB image of the Mississippi Gulfport dataset with selected ROI and (b) Ground truth materials in the ROI with their mean spectra.

 $\label{table V} \textbf{Abundance and Endmember Errors for the Gulfport Dataset}$

	$\times 10^{-4}$	GMM	NCM	NCM w/o PCA	BCM
j	Asphalt	205 \ 52a	1693 \ 94	939 \ 59	1420
	Grass	169 \ 58	1982 \ 121	558 \ 65	2145
	Shadow	499 \ 49	1294 \ 68	921 \ 43	1315
Ì	Tree	1029 \ 89	2194 \ 234	1106 \ 185	2279
_	Roof	908 \ 76	2143 \ 174	1234 \ 104	1657
	Mean	562 \ 65	1861 \ 138	952 \ 91	1763

^a the numbers in ".\." denote the abundance and endmember errors.

less errors than GMM, the difference is less than 0.005 hence negligible.

2) Unsupervised: We created two synthetic images in this case, the first was used to validate the ability to estimate the distribution parameters on scenes with regions of pure pixels, the second was used to validate the segmentation strategy on images with insufficient pure pixels. They were both of size 60×60 pixels and constructed from 4 endmember classes: limestone, basalt, concrete, asphalt, whose spectral signatures were highly differentiable. We assumed that the endmembers were sampled from GMMs following the example in Section II-C. The means of the GMMs were from the ASTER spectral library [50] (see Fig. 4(c) for their spectra) with slight constant changes, which determined a spectral range from $0.4~\mu m$ to $14~\mu m$, re-sampled into 200 values. The covariance matrices were constructed by $a_{jk}^2 \mathbf{I}_B + b_{jk}^2 \mathbf{u}_{jk} \mathbf{u}_{jk}^T$ where \mathbf{u}_{jk} was a unit vector controlling the major variation direction. For the first image, we assumed the 4 materials occupied the 4 quadrants of the square image as pure pixels. Then Gaussian smoothing was applied on each abundance map to make the boundary pixels of each quadrant be mixed by the neighboring materials. For the second image, we made the first material as background, the other materials randomly placed on this background. The procedure of generating the abundance

maps followed [37]: for each material (not as background), 150 Gaussian blobs were randomly placed, whose location and shape width were both sampled from Gaussian distributions. Finally, noise produced similar to above with $\sigma_Y = 0.001$ was added to the generated pixels. Fig. 4 shows the abundance maps, the original spectra of these materials, and the resulting color images by extracting the bands corresponding to wavelengths 488 nm, 556 nm, 693 nm.

The parameters of GMM were $r_{se} = 5$ for the two images, $\beta_1 = 0.1$, $\beta_2 = 0.1$ for the second image. Fig. 5 shows the histograms of ground truth pure pixels and the estimated distributions for the first image. The ground truth distribution is barely visible as most of the time it coincides with GMM. For limestone and asphalt, all the distributions are similar since the pure pixels are generated by a unimodal Gaussian. However, for basalt and concrete, GMM provides a more accurate estimation while the two NCMs seem inferior due to the single Gaussian assumption. The quantitative analysis in Table II implies a similar result by calculating the L_2 distance between the estimated distribution and the ground truth.

Table III shows the comparison of abundance errors from the two images. Since the second image is much more challenging than the first one, we can expect increased errors from all the methods. In general, the results of BCM and the two NCMs show slightly inferior abundances compared to GMM despite the fact that they have access to pure pixels in the image to train their models.

B. Pavia University

The Pavia University dataset was recorded by the Reflective Optics System Imaging Spectrometer (ROSIS) during a flight over Pavia, northern Italy. The dimension is 340 by 610 with a spatial resolution of 1.3 meters/pixel. It has 103 bands with wavelengths ranging from 430 nm to 860 nm. As Fig. 1 shows, the original image contains several man-made and natural materials. Considering that the whole dataset contains many different objects, we only performed experiments on the exemplar ROI (47 by 106) shown in Fig. 1, in which 5 endmembers, meadows, bare soil, painted metal sheets, shadows and pavement, are manually identified.

The parameter of GMM was $r_{se} = 2$. Fig. 6 shows the GMM in the wavelength-reflectance space, where we can see the centers and the major variations of the Gaussians. Fig. 7 shows the scatter plot of the results in the projected space. The scatter plot shows that the identified Gaussian components cover the ground truth pure pixels very well. For painted metal sheets, which has a broad range of pure pixels, it estimated 4 components to cover them. For shadows, only one component was estimated. Fig. 8 shows the histograms of pure pixels and the estimated distributions of GMM and NCMs. We can see that GMM matches the background histogram better than NCMs.

Fig. 9 shows the abundance map comparison. Comparing them with the ground truth shown in Fig. 1(a), we can see that BCM failed to estimate the pure pixels of painted metal sheets, although ground truth pure pixels were used for training.

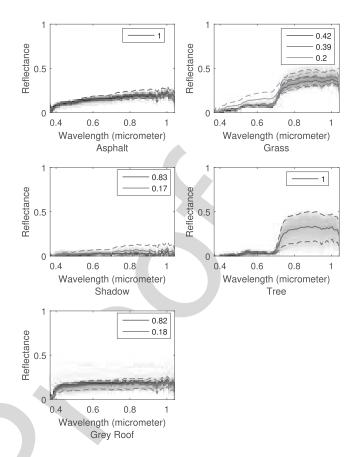


Fig. 11. Estimated GMM in the wavelength-reflectance space for the Mississippi Gulfport dataset. The background gray image and the curves have the same meaning as in Fig. 6.

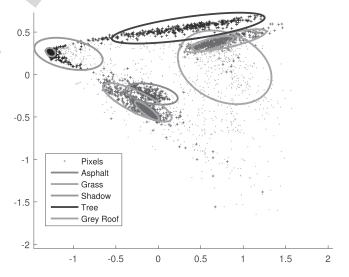


Fig. 12. Scatter plot of the Mississippi Gulfport dataset with the estimated GMM. The ellipses and the dots have the same meaning as in Fig. 7.

For example, the third and fourth abundance maps of BCM show that the pixels in the lower part of painted metal sheets are mixed with shadows, while the reduced reflectances are only caused by angle variation. The result of GMM not only shows sparse abundances for that region, but also interprets the boundary as a combination of neighboring materials. Since this

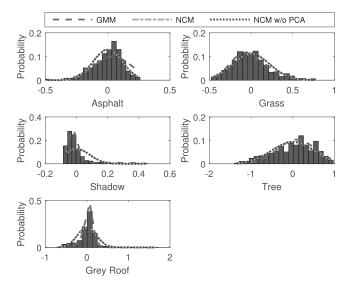


Fig. 13. Histograms of pure pixels for the Gulfport dataset and the estimated distributions from GMM and NCM when projected to 1 dimension.

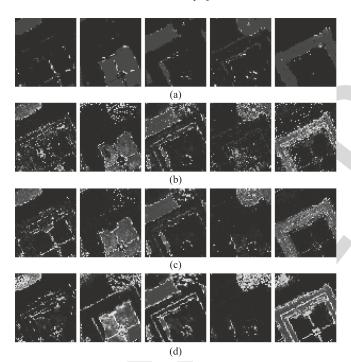


Fig. 14. Abundance maps for the Gulfport dataset. The corresponding endmembers from left to right are asphalt, grass, shadow, tree and grey roof. (a) GMM. (b) NCM. (c) NCM w/o PCA. (d) BCM.

dataset has a spatial spacing of 1.3 meters/pixel, we think this soft transition is more realistic than a simple segmentation. Although the results of NCMs look good in general, the abundances in a pure material region are inconsistent. The errors of abundances and endmembers for these algorithms are shown in Table IV, which implies that GMM performed best overall.

C. Mississippi Gulfport

The dataset was collected over the University of Southern Mississippis-Gulfpark Campus [51]. It is a 271 by 284 image with 72 bands corresponding to wavelengths 0.368 μm to 1.043 μm . The spatial resolution is 1 meter/pixel. The scene

contains several man-made and natural materials including sidewalks, roads, various types of building roofs, concrete, shrubs, trees, and grasses. Since the scene contains many cloths for target detection, we tried to avoid the cloths and selected a 58 by 65 ROI that contains 5 materials [52]. The original RGB image and the selected ROI are shown in Fig. 10(a) while the identified materials and the mean spectra are shown in (b).

The parameter of GMM was $r_{se} = 1$. Fig. 11 shows the GMM result in the wavelength-reflectance space and Fig. 12 shows the scatter plot. We can see that the estimated Gaussian components successfully cover the identified pure pixels. Fig. 13 shows the estimated distributions. Although there are no multiple peaks in any of the histograms, NCMs still do not fit the histograms of shadow and gray roof. In contrast, GMM gives a much better fit for these 2 endmember distributions.

Fig. 14 shows the abundance maps from different algorithms. We can see that GMM matches the ground truth in Fig. 10(b) best, followed by NCM without PCA. This is also verified in the quantitative analysis in Table V. Although NCM and BCM take ground truth pure pixels as input, the scattered dots for trees (fourth abundance map) in both of them and the incomplete region of grass for NCM (asphalt for BCM) show their insufficiency in this case.

V. DISCUSSION AND CONCLUSION

In this paper, we introduced a GMM approach to represent endmember variability, by observing that the identified pure pixels in real applications usually can not be well fitted by a unimodal distribution as in NCM or BCM. We solved several obstacles in linear unmixing using this distribution, including (i) deriving the conditional probability density function of the mixed pixel given each endmember modeled as GMM from two perspectives; (ii) estimating the abundances and endmember distributions by maximizing the log-likelihood with a prior enforcing abundance smoothness and sparsity; (iii) estimating the endmembers for each pixel given the abundances and distribution parameters. The results on synthetic and real datasets show superior accuracy compared to current popular methods like NCM, BCM. Here we have some final remarks.

A. Complexity

As analyzed in Section III-F, each iteration in the estimation of abundances has spatial complexity $O\left(|\mathcal{K}|\,NB^2\right)$ and time complexity $O\left(|\mathcal{K}|\,NB^3\right)$. For comparison, the implemented NCM has the same complexity but with $|\mathcal{K}|=1$. For the supervised synthetic dataset which contains 60 images, the total running time of GMM was 9709 seconds, on a desktop with a Intel Core i7-3820 CPU and 64 GB memory. For comparison, the running time of NCM, NCM without PCA, and BCM was 941, 50751, 62525 seconds respectively. In real applications, running GMM on the Pavia University and Mississippi Gulfport ROIs required 734 seconds and 97 seconds respectively for abundance estimation (24 seconds and 17 seconds for endmember estimation), compared to 40 and 34 seconds from NCM, 1389 and 396 seconds from

903

904

905

906

907

908

909 910

911

912

913

915

917

918

919

921

922

923

924

925

926

927

929

931

932

933

934

935

937

938

939

941

942

943

945

946

947

948

949

950

951

952

953

954

956

958

960

962

963

964

966

967

968

969

970

971

973

974

975

851

853

855

856

857

858

860

862

864

866

868

870

871

872

874

875

876

878

880

882

884

886

890

892

893

895

896

897

898

899

900

NCM without PCA, 1170 and 616 seconds from BCM. As analyzed, the main factors affecting the efficiency of GMM and NCMs are $|\mathcal{K}|$ and B.

B. Limitation

The complexity analysis leads to one limitation of the method. That is, the complexity grows exponentially with increasing numbers of components. This could cause problems for a large amount of pure pixels. To overcome this shortcoming, there are some empirical workarounds, such as reducing the number of components by introducing thresholds, or reducing the number of pure pixels to a fixed number by random sampling. Another limitation is that the proposed unsupervised version assumes presence of regions of pure pixels, which mostly happens in urban scenes. For scenes with a lot of mixed pixels, this assumption may not hold. Note that unsupervised unmixing is a very challenging problem. The previous works for this problem all assume several properties on the abundances and endmembers [21]-[23]. Hence, this limitation exists more or less in all the works on this problem. Finally, the method was only evaluated on real urban datasets with only ground truth on pure pixels: it is therefore unclear if the abundance estimation on mixed pixels is also accurate. This is due to lack of datasets and ground truth in the hyperspectral community. We plan to validate it on a more comprehensive dataset given in [31] in the future.

C. Future Work

The proposed GMM formulation has several applications that we can investigate in the future. First, in target detection, endmember variability may interfere with the target as well as the background [53]. By modeling the target or the background as spectra sampled from GMM distributions, we may devise more sophisticated and accurate target detection algorithms. Second, in fusion of hyperspectral and multispectral images, the LMM is usually used to overcome the underdetermined nature of the problem [54], [55]. However, the LMM does not hold in real scenarios as shown in this work. If we use the LMM with endmember variability, which is modeled by samples from GMM distributions, we may have a fusion algorithm that better fits the data. Finally, in estimating the noise or intrinsic dimension of hyperspectral images, simulated data are generated to quantify the results [46]. When these simulated data are created, usually the LMM is used without considering the endmember variability. Using the GMM formulation, we may generate distinct endmembers for each pixel and create more realistic synthetic data.

REFERENCES

- [1] M. Berman, H. Kiiveri, R. Lagerstrom, A. Ernst, R. Dunne, and J. F. Huntington, "ICE: A statistical approach to identifying endmembers in hyperspectral images," *IEEE Trans. Geosci. Remote Sens.*, vol. 42, no. 10, pp. 2085–2095, Oct. 2004.
- [2] J. M. P. Nascimento and J. M. Bioucas-Dias, "Vertex component analysis: A fast algorithm to unmix hyperspectral data," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 4, pp. 898–910, Apr. 2005.
- [3] A. Zare, P. Gader, O. Bchir, and H. Frigui, "Piecewise convex multiple-model endmember detection and spectral unmixing," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 5, pp. 2853–2862, May 2013.

- [4] J. M. Bioucas-Dias et al., "Hyperspectral unmixing overview: Geometrical, statistical, and sparse regression-based approaches," IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens., vol. 5, no. 2, pp. 354–379, Apr. 2012.
- [5] N. Keshava and J. F. Mustard, "Spectral unmixing," IEEE Signal Process. Mag., vol. 19, no. 1, pp. 44–57, Jan. 2002.
- [6] B. Hapke, "Bidirectional reflectance spectroscopy: 1. Theory," J. Geophys. Res., Solid Earth, vol. 86, no. B4, pp. 3039–3054, 1981
- [7] A. Halimi, Y. Altmann, N. Dobigeon, and J.-Y. Tourneret, "Nonlinear unmixing of hyperspectral images using a generalized bilinear model," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 11, pp. 4153–4162, Nov. 2011.
- [8] B. Somers et al., "Nonlinear hyperspectral mixture analysis for tree cover estimates in orchards," *Remote Sens. Environ.*, vol. 113, no. 6, pp. 1183–1193, Feb. 2009.
- [9] R. Heylen and P. D. Gader, "Nonlinear spectral unmixing with a linear mixture of intimate mixtures model," *IEEE Geosci. Remote Sens. Lett.*, vol. 11, no. 7, pp. 1195–1199, Jul. 2014.
- [10] J. Broadwater and A. Banerjee, "A generalized kernel for areal and intimate mixtures," in *Proc. 2nd Workshop Hyperspectral Image Signal Process., Evol. Remote Sens. (WHISPERS)*, Jun. 2010, pp. 1–4.
- [11] J. Broadwater, R. Chellappa, A. Banerjee, and P. Burlina, "Kernel fully constrained least squares abundance estimates," in *Proc. IEEE Int. Geosci. Remote Sens. Symp. (IGARSS)*, Jul. 2007, pp. 4041–4044.
- [12] R. Heylen, M. Parente, and P. Gader, "A review of nonlinear hyperspectral unmixing methods," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 6, pp. 1844–1868, Jun. 2014.
- [13] B. Somers, G. P. Asner, L. Tits, and P. Coppin, "Endmember variability in spectral mixture analysis: A review," *Remote Sens. Environ.*, vol. 115, no. 7, pp. 1603–1616, 2011.
- [14] A. Zare and K. Ho, "Endmember variability in hyperspectral analysis: Addressing spectral variability during spectral unmixing," *IEEE Signal Process. Mag.*, vol. 31, no. 1, pp. 95–104, Jan. 2014.
- [15] X. Du, A. Zare, P. Gader, and D. Dranishnikov, "Spatial and spectral unmixing using the beta compositional model," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 6, pp. 1994–2003, Jun. 2014.
- [16] A. Zare and P. Gader, "PCE: Piecewise convex endmember detection," IEEE Trans. Geosci. Remote Sens., vol. 48, no. 6, pp. 2620–2632, Jun. 2010.
- [17] D. A. Roberts, M. Gardner, R. Church, S. Ustin, G. Scheer, and R. O. Green, "Mapping chaparral in the Santa Monica Mountains using multiple endmember spectral mixture models," *Remote Sens. Environ.*, vol. 65, no. 3, pp. 267–279, Sep. 1998.
- [18] A. Halimi, N. Dobigeon, and J.-Y. Tourneret, "Unsupervised unmixing of hyperspectral images accounting for endmember variability," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 4904–4917, Dec. 2015.
- [19] O. Eches, N. Dobigeon, C. Mailhes, and J.-Y. Tourneret, "Bayesian estimation of linear mixtures using the normal compositional model. Application to hyperspectral imagery," *IEEE Trans. Image Process.*, vol. 19, no. 6, pp. 1403–1413, Jun. 2010.
- [20] C. A. Bateson, G. P. Asner, and C. A. Wessman, "Endmember bundles: A new approach to incorporating endmember variability into spectral mixture analysis," *IEEE Trans. Geosci. Remote Sens.*, vol. 38, no. 2, pp. 1083–1094. Mar. 2000.
- [21] L. Drumetz, M.-A. Veganzones, S. Henrot, R. Phlypo, J. Chanussot, and C. Jutten, "Blind hyperspectral unmixing using an extended linear mixing model to address spectral variability," *IEEE Trans. Image Process.*, vol. 25, no. 8, pp. 3890–3905, Aug. 2016.
- [22] P.-A. Thouvenin, N. Dobigeon, and J.-Y. Tourneret, "Hyperspectral unmixing with spectral variability using a perturbed linear mixing model," *IEEE Trans. Signal Process.*, vol. 64, no. 2, pp. 525–538, Jan. 2016.
- [23] A. Halimi, P. Honeine, and J. M. Bioucas-Dias, "Hyperspectral unmixing in presence of endmember variability, nonlinearity, or mismodeling effects," *IEEE Trans. Image Process.*, vol. 25, no. 10, pp. 4565–4579, Oct. 2016.
- [24] B. Zhang, L. Zhuang, L. Gao, W. Luo, Q. Ran, and Q. Du, "PSO-EM: A hyperspectral unmixing algorithm based on normal compositional model," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 12, pp. 7782–7792, Dec. 2014.
- [25] O. Eches, N. Dobigeon, and J.-Y. Tourneret, "Estimating the number of endmembers in hyperspectral images using the normal compositional model and a hierarchical Bayesian algorithm," *IEEE J. Sel. Topics Signal Process.*, vol. 4, no. 3, pp. 582–591, Jun. 2010.

1056

1057

1059

1060

1061

1063

1064

1065

1066

1068

1069

1070

1071

1073

1074

1075

1077

1078

1079

1080

1082

1083

1085

1087

1089

1090

1092

1093

1094

1095

1097

1099

1100

1102

1104

1105

1106

1107

1108

1058 AO:3

- [26] C. Song, "Spectral mixture analysis for subpixel vegetation fractions in the urban environment: How to incorporate endmember variability?" *Remote Sens. Environ.*, vol. 95, no. 2, pp. 248–263, 2005.
- [27] J.-P. Combe et al., "Analysis of OMEGA/Mars express data hyper-spectral data using a multiple-endmember linear spectral unmixing model (MELSUM): methodology and first results," Planetary Space Sci., vol. 56, no. 7, pp. 951–975, May 2008.
- [28] G. P. Asner and D. B. Lobell, "A biogeophysical approach for automated SWIR unmixing of soils and vegetation," *Remote Sens. Environ.*, vol. 74, no. 1, pp. 99–112, Oct. 2000.
- [29] G. P. Asner and K. B. Heidebrecht, "Spectral unmixing of vegetation, soil and dry carbon cover in arid regions: Comparing multispectral and hyperspectral observations," *Int. J. Remote Sens.*, vol. 23, no. 19, pp. 3939–3958, Oct. 2002.
- [30] A. Castrodad, Z. Xing, J. B. Greer, E. Bosch, L. Carin, and G. Sapiro, "Learning discriminative sparse representations for modeling, source separation, and mapping of hyperspectral imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 11, pp. 4263–4281, Nov. 2011.
- [31] E. B. Wetherley, D. A. Roberts, and J. P. McFadden, "Mapping spectrally similar urban materials at sub-pixel scales," *Remote Sens. Environ.*, vol. 195, pp. 170–183, Jun. 2017.
- [32] D. Stein, "Application of the normal compositional model to the analysis of hyperspectral imagery," in *Proc. IEEE Workshop Adv. Techn. Anal. Remotely Sensed Data*, Oct. 2003, pp. 44–51.
- [33] L. Tits, B. Somers, and P. Coppin, "The potential and limitations of a clustering approach for the improved efficiency of multiple endmember spectral mixture analysis in plant production system monitoring," *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 6, pp. 2273–2286, Jun. 2012.
- [34] M.-D. Iordache, L. Tits, J. M. Bioucas-Dias, A. Plaza, and B. Somers, "A dynamic unmixing framework for plant production system monitoring," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 6, pp. 2016–2034, Jun. 2014.
- [35] L. Tits, B. Somers, W. Saeys, and P. Coppin, "Site-specific plant condition monitoring through hyperspectral alternating least squares unmixing," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 8, pp. 3606–3618, Aug. 2014.
- [36] J. B. Lee, A. S. Woodyatt, and M. Berman, "Enhancement of high spectral resolution remote-sensing data by a noise-adjusted principal components transform," *IEEE Trans. Geosci. Remote Sens.*, vol. 28, no. 3, pp. 295–304, May 1990.
- [37] Y. Zhou, A. Rangarajan, and P. D. Gader, "A spatial compositional model for linear unmixing and endmember uncertainty estimation," *IEEE Trans. Image Process.*, vol. 25, no. 12, pp. 5987–6002, Dec. 2016.
- [38] D. Achlioptas and F. McSherry, "On spectral learning of mixtures of distributions," in *Learning Theory*. Springer, 2005, pp. 458–469.
- [39] N. Vlassis and A. Likas, "A greedy EM algorithm for Gaussian mixture learning," *Neural Process. Lett.*, vol. 15, no. 1, pp. 77–87, Feb. 2002.
- [40] K. Lange, Optimization. Springer, 2013.
- [41] X.-L. Meng and D. B. Rubin, "Maximum likelihood estimation via the ECM algorithm: A general framework," *Biometrika*, vol. 80, no. 2, pp. 267–278, 1993.
- [42] D. P. Bertsekas, Nonlinear Programming. Belmont, MA, USA: Athena Scientific, 1999.
- [43] C.-J. Lin, "Projected gradient methods for nonnegative matrix factorization," Neural Comput., vol. 19, no. 10, pp. 2756–2779, 2007
- [44] G. J. McLachlan and S. Rathnayake, "On the number of components in a Gaussian mixture model," Wiley Interdiscipl. Rev., Data Mining Knowl. Discovery, vol. 4, no. 5, pp. 341–355, 2014.
- [45] P. Smyth, "Model selection for probabilistic clustering using cross-validated likelihood," *Statist. Comput.*, vol. 10, no. 1, pp. 63–72, Jan. 2000.
- [46] L. Gao, Q. Du, B. Zhang, W. Yang, and Y. Wu, "A comparative study on linear regression-based noise estimation for hyperspectral imagery," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 6, no. 2, pp. 488–498, Apr. 2013.
- [47] R. E. Roger, "Principal Components transform with simple, automatic noise adjustment," *Int. J. Remote Sens.*, vol. 17, no. 14, pp. 2719–2727, 1996.
- [48] N. Guan, D. Tao, Z. Luo, and B. Yuan, "Manifold regularized discriminative nonnegative matrix factorization with fast gradient descent," *IEEE Trans. Image Process.*, vol. 20, no. 7, pp. 2030–2048, Jul. 2011.
- [49] C. Bishop, Pattern Recognition and Machine Learning. New York, NY, USA: Springer, 2006.

- [50] A. M. Baldridge, S. J. Hook, C. I. Grove, and G. Rivera, "The ASTER spectral library version 2.0," *Remote Sens. Environ.*, vol. 113, no. 4, pp. 711–715, 2009.
- [51] P. Gader, A. Zare, R. Close, J. Aitken, and G. Tuell, "MUUFL Gulfport hyperspectral and LiDAR airborne data set," Univ. Florida, Gainesville, FL, USA, Tech. Rep. REP-2013-570, 2013.
- [52] X. Du and A. Zare, "Technical report: Scene label ground truth map for MUUFL Gulfport data set," Univ. Florida, Gainesville, FL, USA, Tech. Rep. 20170417, 2017.
- [53] C. Jiao and A. Zare, "Functions of multiple instances for learning target signatures," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 8, pp. 4670–4686, Aug. 2015.
- [54] N. Yokoya, T. Yairi, and A. Iwasaki, "Coupled nonnegative matrix factorization unmixing for hyperspectral and multispectral data fusion," *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 2, pp. 528–537, Feb. 2012.
- [55] Q. Wei, N. Dobigeon, and J.-Y. Tourneret, "Fast fusion of multi-band images based on solving a Sylvester equation," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 4109–4121, Nov. 2015.



Yuan Zhou received the B.E degree in software engineering and the M.E. degree in computer application technology from the Huazhong University of Science and Technology, Wuhan, China, in 2008 and 2011, respectively. Since 2013, he is currently pursuing the Ph.D. degree with the Department of CISE, University of Florida, Gainesville, FL, USA. He was with Shanghai UIH as a Software Engineer for two years. His research interests include image processing, computer vision, and machine learning.



Anand Rangarajan is currently with the Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL, USA. His research interests are machine learning, computer vision, and the scientific study of consciousness.



Paul D. Gader (M'86–SM'09–F'11) received the Ph.D. degree in mathematics for image-processing-related research from the University of Florida, Gainesville, FL, USA, in 1986. He was a Senior Research Scientist with Honeywell, a Research Engineer and a Manager with the Environmental Research Institute of Michigan, Ann Arbor, MI, USA, and a Faculty Member with the University of Wisconsin, Oshkosh, WI, USA, the University of Missouri, Columbia, MO, USA, and the University of Florida, FL, USA, where he is currently a

Professor of Computer and Information Science and Engineering. He was a Summer Student Fellow at the Eglin Air Force Base involving in algorithms for the detection of bridges in forward-looking infrared imagery, where he performed his first research in image processing in 1984. He has been involved in a wide variety of theoretical and applied research problems including fast computing with linear algebra, mathematical morphology, fuzzy sets, Bayesian methods, handwriting recognition, automatic target recognition, biomedical image analysis, landmine detection, human geography, and hyperspectral and light detection, and ranging image analysis projects. He has authored or co-authored hundreds of refereed journal and conference papers.

977

978

979

980 981

982

984

985

986

987

988

989

990

991

992

993

994

995

996 997

998

999

1000

1001

1002

1003

1005

1006

1007

1008

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030 1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

AQ:2

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES

PLEASE NOTE: We cannot accept new source files as corrections for your paper. If possible, please annotate the PDF proof we have sent you with your corrections and upload it via the Author Gateway. Alternatively, you may send us your corrections in list format. You may also upload revised graphics via the Author Gateway.

AQ:1 = Please note that there were discrepancies between the accepted pdf [gmm_var_journal_ieee_final.pdf] and the [gmm_var_journal_ieee_final.tex] in the sentence on lines 188–195, 623–632, and references. We have followed [gmm_var_journal_ieee_final.tex].

AQ:2 = Please provide the publisher location for refs. [38] and [40].

AQ:3 = Please provide the department name for refs. [51] and [52].