



Faculty of Sciences
Department of Applied Mathematics, Computer Science and Statistics
Faculty of Engineering
Department of Telecommunications and Information Processing (TELIN)
Image Processing and Interpretation (IPI)

Volume-based geometric and Bayesian approaches in linear hyperspectral unmixing

Davor JOSIPOVIC

Promoter: Prof. Dr. Aleksandra PIZURICA
Co-promoter: Prof. Dr. Hongyan ZHANG
Co-promoter: Prof. Dr. Dries BENOIT

Thesis submitted in fulfillment of the requirements for the degree of
MASTER OF SCIENCE IN STATISTICAL DATA ANALYSIS

Year: 2016–2017

Compiled: June 21, 2017



Faculty of Sciences
Department of Applied Mathematics, Computer Science and Statistics
Faculty of Engineering
Department of Telecommunications and Information Processing (TELIN)
Image Processing and Interpretation (IPI)

Volume-based geometric and Bayesian approaches in linear hyperspectral unmixing

Davor JOSIPOVIC

Promoter: Prof. Dr. Aleksandra PIZURICA
Co-promoter: Prof. Dr. Hongyan ZHANG
Co-promoter: Prof. Dr. Dries BENOIT

Thesis submitted in fulfillment of the requirements for the degree of
MASTER OF SCIENCE IN STATISTICAL DATA ANALYSIS

Year: 2016–2017

Compiled: June 21, 2017

The author and the promoter give permission to consult this master dissertation and to copy it or parts of it for personal use. Each other use falls under the restrictions of the copyright, in particular concerning the obligation to mention explicitly the source when using results of this master dissertation.

Abstract

Hyperspectral cameras make images in hundreds or thousands of spectral channels. Due to low spatial resolution, spectra measured by such cameras are mixtures of spectra of materials in the scene. The reverse process in which materials (called endmembers) are estimated is called hyperspectral unmixing. Unmixing is challenging due to model underspecification, observation noise, endmember variability and high dimensionality. The main focus of this thesis is understanding the statistical assumptions, derivations and implementation details of the geometrical and Bayesian approaches for unmixing. For the geometrical approach we mainly focus on Berman's ICE algorithm. For the Bayesian approach we focus on Arngren's BayesNMF-Vol algorithm. We start by defining the linear mixing model. Then we derive the maximum likelihood estimators for the endmembers and abundances. We show the equivalence between maximizing this log-likelihood function and minimizing the objective function of the ICE algorithm. We further analyze every detail of the minimization process of this function. The limitations drive us to more complex approaches. That is why we subsequently put the unmixing problem into a Bayesian framework. Solving for the posterior calls for custom made Gibbs samplers which we cover in great detail. Finally, the algorithm characteristics are illustrated using synthetic and real data.

The most important original contributions of this work are: (1) the in-depth elaborate analysis of Berman's ICE and Arngren's BayesNMF-Vol algorithm, with original derivations and insights that promote understanding of these methods and their limitations; (2) a novel unmixing algorithm, dubbed ICE-S, that extends the ICE algorithm with spatial information and (3) an original R-framework for the analysis of hyperspectral images, build from ground up and made available as open source code.

Preface

I started working on a topic of deep learning in image pixel classification mainly as a way to learn more about machine learning approaches that I found lacking in the excellent Master of Statistical Data Analysis programme at the University of Ghent. The purpose was to associate each pixel of a hyperspectral image with a class, for example water, brick, grass, etc. Only about 5 % of the image was labeled by a human. The rest of the pixels were for the algorithm to label. The deep learning approaches were mainly focused on prediction, and as such, 'anything goes'. I understood why a logistic regression classifier was a good fit, and even a 1-layer neural network – because it allowed for multiple different pixel-spectra to be associated with the same class – but the state-of-the art approaches did lots of very strange things that apparently worked better than these classic approaches. But why? Were they modeling human error? A pixel in the middle of an ocean could be composed of plastic too, like from a small buoy. If I was able to somehow decompose this pixel in plastic and water, wouldn't that make a better classifier? Little did I know that behind this decomposition there was more than 25 years of research called hyperspectral unmixing. And that is how this thesis was born.

I wish to thank my promoter Prof. Aleksandra Pizurica for the constant encouragements, suggestions, follow up and multiple readthroughs, Prof. Hongyan Zhang for introducing me to the problem of hyperspectral pixel classification and Prof. Dries Benoit for explaining me conceptually the reason behind the intrinsic regularization of the Bayesian approach.

Contents

1	Introduction	7
1.1	Hyperspectral imaging	8
1.2	Spectral mixing	8
1.2.1	Linear mixing model	9
1.2.2	Nonlinear mixing model	11
1.3	Hyperspectral unmixing process	11
1.3.1	Unmixing methods	12
1.3.2	Spatial information	13
1.4	Applications	13
1.5	Notation and conventions	13
1.6	Thesis organization and contributions	14
2	ICE algorithm and extensions	15
2.1	Objective function	16
2.2	Regularized objective function	17
2.3	Minimization	18
2.3.1	Coordinate descent	18
2.3.2	Block coordinate descent	19
2.3.3	Minimize $L_{reg}(\mathbf{w}_n)$ with fixed \mathbf{E} and $\mathbf{W}_{\tilde{n}}$:	19
2.3.4	Minimize $L_{reg}(\mathbf{e}_b)$ with fixed \mathbf{W} and $\mathbf{E}_{\tilde{b}}$:	21
2.3.5	The ICE algorithm	22
2.3.6	Non-Convexity of L_{reg}	22
2.4	Model peculiarities	23
2.4.1	L dependence on B	23
2.4.2	Disregard of σ^2	24
2.5	Extension with spatial regulation (ICE-S)	24
2.5.1	Motivation	24
2.5.2	New regularization function S	25
2.5.3	New objective function L_{VS}	25
2.5.4	Minimization	26
2.5.5	Derivation of $L_{VS}(\mathbf{E})$	26
2.6	Evaluation	26
2.7	Conclusion	27
3	BayesNMF-Vol algorithm	29
3.1	Belief Network for the Bayesian model	29
3.2	The partial probability density functions	30
3.2.1	The likelihood $f_{\mathbf{X} \mathbf{E},\mathbf{W},\sigma^2}$	31

3.2.2	The endmember prior $f_{\mathbf{E} \gamma}$	31
3.2.3	The abundances prior $f_{\mathbf{W}}$	31
3.2.4	The variance prior $f_{\sigma^2 \alpha,\beta}$	31
3.2.5	The posterior $f_{\mathbf{E},\mathbf{W},\sigma^2 \mathbf{X},\alpha,\beta,\gamma}$	32
3.3	Relation between MAP and ICE estimators	33
3.4	The BayesNMF-Vol sampler algorithm	33
3.4.1	Choice of a sampler	33
3.4.2	The Gibbs sampler	33
3.4.3	Sampling σ^2	34
3.4.4	Sampling \mathbf{W}	35
3.4.5	Sampling \mathbf{E}	36
3.5	Model peculiarities	37
3.5.1	Intrinsic regularization of the Bayesian approach	37
3.6	Conclusion	39
4	Synthetic data	40
4.1	Image generation procedure	40
4.2	Analysis	41
4.3	Evaluation	41
4.3.1	Error metric selection	41
4.3.2	Results	42
4.4	Conclusion	43
5	Real AVIRIS data - Cuprite	45
5.1	Preprocessing AVIRIS data	45
5.1.1	Raw data	45
5.1.2	Radiance data	45
5.1.3	Reflectance data	46
5.1.4	Data cleanup	47
5.2	Scene selection	48
5.3	Analysis	48
5.4	Results	49
5.5	Conclusion	50
6	Conclusion	53
	Appendices	55
A	ICE appendix	56
A.1	Relation between V and SSD	56
A.2	Minimum of $L_{reg}(\mathbf{E})$ given \mathbf{W}	57
A.3	Derivation of $S_n(\mathbf{w}_n)$	58
B	BayesNMF-Vol appendix	59
B.1	Density of linearly transformed variable	59
B.2	Derivation of $\mathcal{N}(x_k \mathbf{x}_{\tilde{k}}, \hat{\mu}_k, \hat{\sigma}_k^2)$ from $\mathcal{N}(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma})$	59
B.3	Derivation of $f_{\mathbf{W} \mathbf{E},\sigma^2,\mathbf{X}}$	60
B.4	Mapping $f_{\mathbf{w}_n \mathbf{E},\sigma^2,\mathbf{x}_n}$ to $f_{\mathbf{y} \mathbf{E},\sigma^2,\mathbf{x}_n}$	61
B.5	Derivation of $f_{y_k \mathbf{y}_{\tilde{k}},\mathbf{E},\sigma^2,\mathbf{x}_n}$ for Gibbs sampling	62

B.6	Derivation of $f_{\mathbf{E} \mathbf{W},\sigma^2,\mathbf{X},\gamma}$	63
B.7	Derivation of $f_{e_{bm} e_{b\bar{m}},\mathbf{W},\sigma^2,\mathbf{x}_b,\gamma}$ for Gibbs sampling	64
C	Source code	66

Chapter 1

Introduction

In this thesis our main focus is *understanding* the methods used in hyperspectral unmixing. This will lead us to some new insights into the current methods, their limitations and some extensions and improvements aimed at overcoming these limitations. We start from the following questions: What is a hyperspectral image? How does it facilitate material detection? What is hyperspectral unmixing? And most importantly, how is it done?

In this Chapter we first learn about the hyperspectral image. Then we look at how materials (that compose the image) leave mixed “fingerprints” in the pixel spectra. Our goal is to extract them. But to be able to extract them, we first need a mathematical model of how they are mixed. Such a mathematical model is defined in Section 1.2. Once we have the mixing model, we explain the unmixing process in Section 1.3 together with some state-of-the-art in the field. Subsequently, in Section 1.4, we give some examples of unmixing applications. In Section 1.5 we cover the notational conventions used in this thesis. And finally in Section 1.6 we give an overview of the subsequent chapters and main contributions.

So what will one learn from these other chapters? Well, if you read the original papers you will note that they are quite succinct on various important topics. For example, when Berman [5] defines the objective function which is to be minimized, he just notes that it follows from the mixing model. We actually derive it and show which assumptions he makes underway. Also, Berman explains the minimization process in a few sentences. But the actual minimization of the constrained quadratic function is much more challenging because the positive-definiteness of the Hessian matrix can not always be guaranteed. Similarly, we find that in the Bayesian model the key underlying assumption is a directed acyclic graph (DAG) with the Markov property without which none of the derivations for the Gibbs sampler can be made. Furthermore, sampling from a degenerate distribution is a challenge in its own. But the original papers put hardly any emphasis on this. So in summary, you will find here all the derivations and nuances which were omitted from the original papers. This will definitely improve the understanding of the discussed unmixing approaches.

Apart from these elaborations and small improvement suggestions, one will also find our own contribution: the extension of the ICE algorithm with spatial information, which we dub ICE-S. But before we delve into the details, we review briefly the basic concepts of hyperspectral imaging.

1.1 Hyperspectral imaging

A hyperspectral image can be seen as an extension of an ordinary RGB image. Whereas an RGB image consists of three bands, red green and blue taken at 485 nm, 550 nm and 645 nm respectively, a hyperspectral image consists of many more (several tens to several hundreds) bands extended beyond the visible spectrum. It is usually represented as a 3D data cube such as the one depicted in Fig. 1.1.

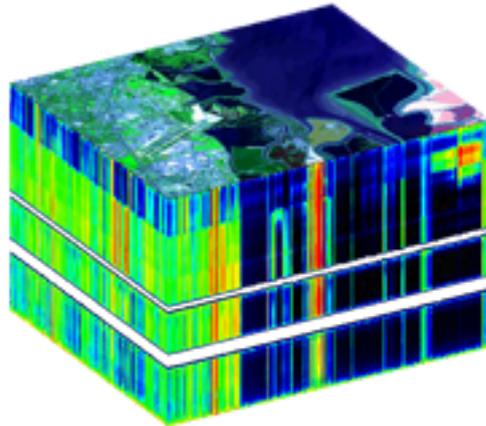


Figure 1.1: This sample data cube [14] was acquired with AVIRIS instrument on August 20, 1992 when it was flown on a NASA ER-2 plane at an altitude of 20,000 meters over Moffett Field, California, at the southern end of the San Francisco Bay. The top of the cube is a false-color image made to accentuate the structure in the water and evaporation ponds on the right. Also visible on the top of the cube is the Moffett Field airport. The sides of the cube are slices showing the edges of the top in all 224 of the AVIRIS spectral bands. The tops of the sides are in the visible part of the spectrum (400 nm), and the bottoms are in the infrared (2500 nm). The sides are pseudo-color, ranging from black and blue (low response) to red (high response).

A hyperspectral image is not to be confused with a multispectral image. For example, the Landsat Thematic Mapper attached to the Landsat 4 and 5 satellites detects only 6 bands (and extra one on a much higher spatial resolution). Researchers use the ratios of those bands to discriminate surface materials with limited accuracy. Hyperspectral images contain many more narrow and continuous bands. This allows better accuracy in analysis. In [11] one can find a more elaborate comparison.

1.2 Spectral mixing

Materials leave unique “fingerprints” in the electromagnetic spectrum, known as spectral signatures. These signatures enable identification of the materials that make up the image. Fig. 1.2 shows this. Notice that there are pure and mixed pixels. For example, the middle emphasized pixel consists entirely of water, while the top emphasized pixel is composed of soil and rock. The latter is due to the spatial resolution of the sensor.

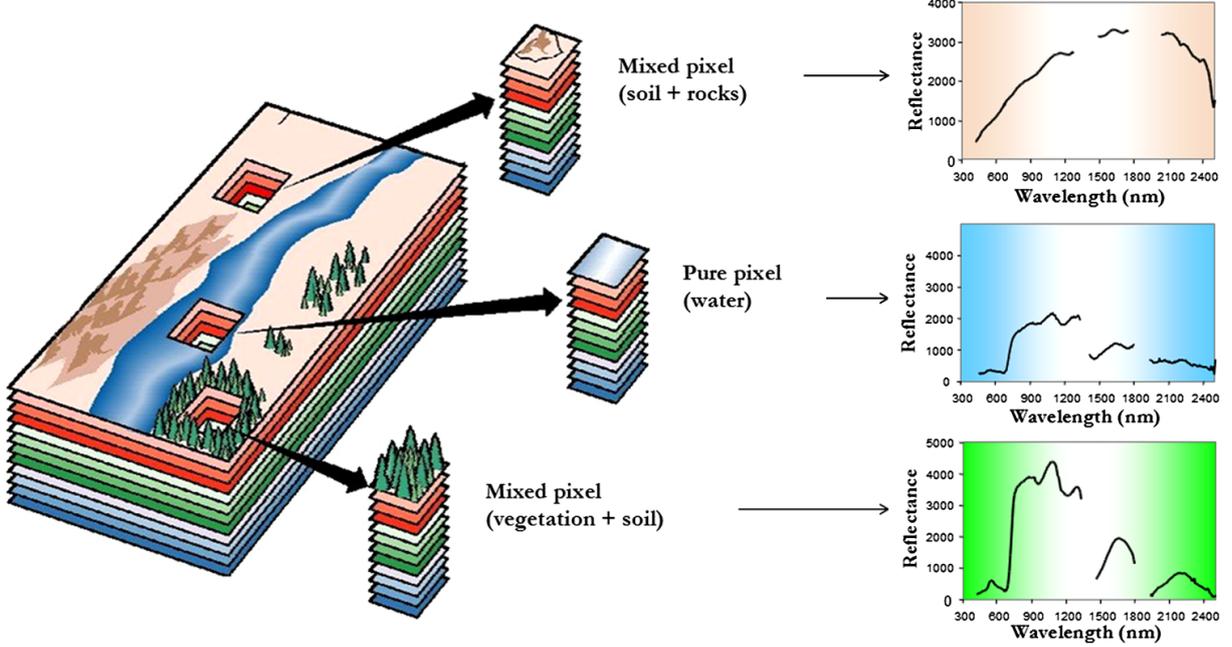


Figure 1.2: Hyperspectral image decomposition [8].

Hyperspectral unmixing is any process that separates the pixel spectra from a hyperspectral image into a collection of constituent spectral signatures, called *endmembers*, and a set of fractional *abundances*, one set per pixel. To relate to Fig. 1.2, the top emphasized pixel is composed of two endmembers (soil and rocks) in abundances of 60 and 40 percent respectively.

1.2.1 Linear mixing model

The linear mixing model (LMM) defines mathematically how the spectra of the materials in a scene are mixed together in a pixel. It assumes that the incident light-rays interact only with the material on which they scatter. In this case, the mixing occurs within the instrument itself due to the fact that the spatial resolution is not fine enough. This mechanism is shown in Fig. 1.3.

Under such a linear mixture, the relative area of the endmembers corresponds to their abundances. Thus the measured spectrum at a pixel $\mathbf{x}_n \in \mathbb{R}_+^B$ is a weighted average of the radiances of the materials (i.e. endmembers \mathbf{e}_m) present at the pixel. We write this as follows:

$$\begin{aligned}
 \mathbf{x}_n &= w_{n1}\mathbf{e}_1 + \dots + w_{nM}\mathbf{e}_M + \boldsymbol{\epsilon} \\
 &= \left(\sum_{m=1}^M w_{nm}\mathbf{e}_m \right) + \boldsymbol{\epsilon} \\
 &= \mathbf{E}\mathbf{w}_n + \boldsymbol{\epsilon}
 \end{aligned} \tag{1.1}$$

where M denotes the number of endmembers in the image; $\mathbf{E} = [\mathbf{e}_1, \dots, \mathbf{e}_M] \in \mathbb{R}_+^{B \times M}$ and $\mathbf{w}_n \in \mathbb{R}_+^{M \times 1}$ is a column vector of endmember proportions in \mathbf{x}_n . The noise $\boldsymbol{\epsilon}$ is usually taken as $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ distributed (i.e. the error terms are assumed independent).

Note also that

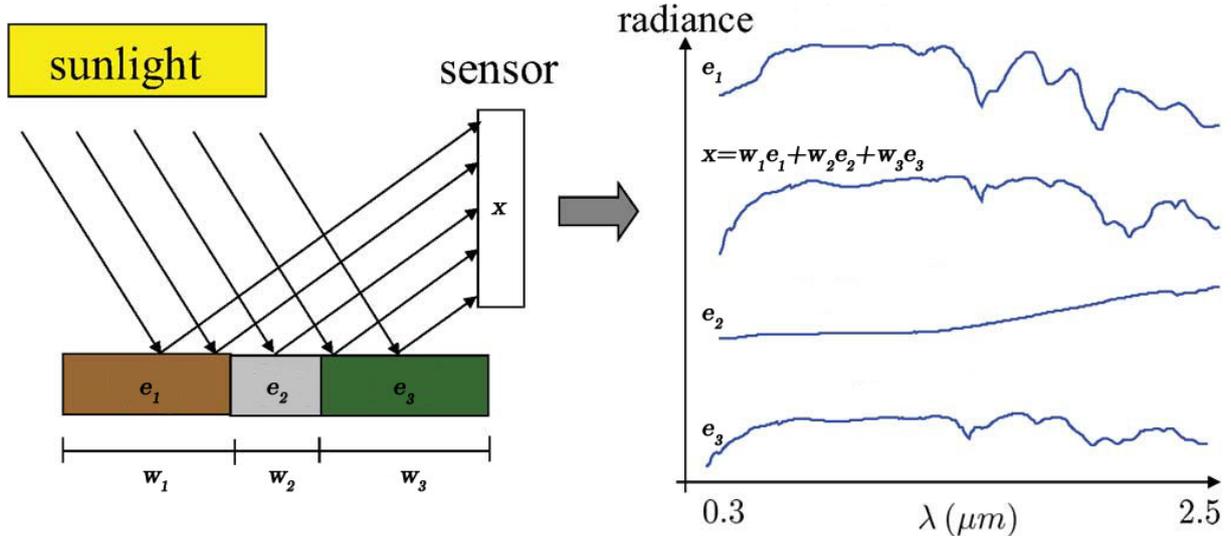


Figure 1.3: Linear mixing mechanics. Sunlight scattered by three materials denoted by $\{e_1, e_2, e_3\}$ in a scene is incident on a sensor-pixel \mathbf{x} that measures radiance in B bands.

$$\sum_{m=1}^M w_{nm} = |\mathbf{w}_n|_1 = 1 \quad (1.2)$$

together with non-negativity of w_{nm} is called the closure or convex geometry constraint [7]. This implies that \mathbf{w}_n lies on a unit $(M - 1)$ -simplex.

Eq. 1.1 can be further generalized to all pixels:

$$\mathbf{X} = \mathbf{W}\mathbf{E}^T + \boldsymbol{\varepsilon} \quad (1.3)$$

where matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T \in \mathbb{R}_+^{N \times B}$ of the observed spectra and matrix $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_N]^T \in \mathbb{R}_+^{N \times M}$ whose n th row represents the endmember concentration profile of n th sample \mathbf{x}_i .

Multidimensional simplex representation

The non-negativity and sum-to-unity constraints in Eq. 1.2 imply that \mathbf{x}_n are contained within a simplex where the endmembers e_m form the vertices. This is shown in Fig 1.4. The fact that some \mathbf{x}_n fall out of the simplex is due to noise. Furthermore, depending on the source, the image may not contain pure pixels (i.e. pixels constituted by only one endmember/material). Fig. 1.4b is such an example.

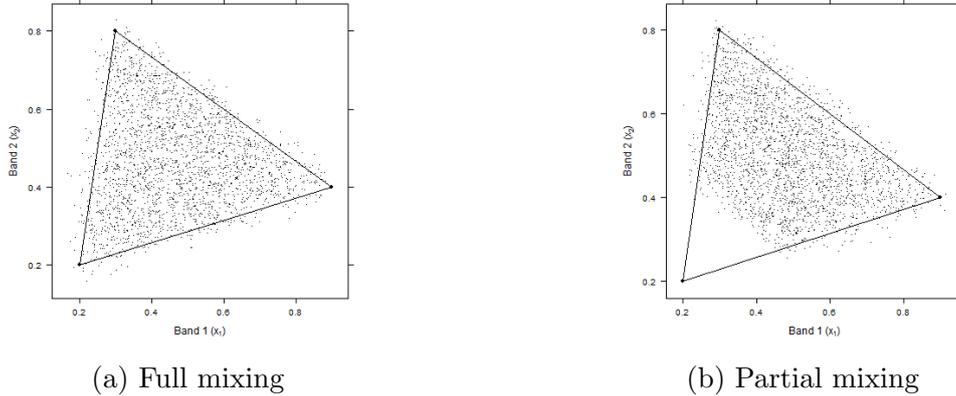


Figure 1.4: 2500 pixels $\mathbf{x}_i \in \mathbb{R}_+^2$ represented in function of their band. The observed pixels may not span the entire simplex due to the mixing of the endmembers. This is shown in (b) where the bottom-left endmember is never available in pure form in the image. Such mixing is called partial.

1.2.2 Nonlinear mixing model

Nonlinear mixing is usually due to physical interactions (i.e. reflections and absorptions) between the light scattered by multiple materials in the scene. We do not consider nonlinear mixing in this thesis. This is partly justified by the fact that linear mixing is an acceptable approximation of the light scattering mechanisms in many real scenarios [8].

1.3 Hyperspectral unmixing process

The hyperspectral unmixing process is a step-by-step procedure of extracting the materials (i.e. endmember spectra) from hyperspectral images. This process usually consists of the following steps:

1. *Atmospheric correction* by which atmospheric effects are compensated by converting radiance to reflectance data. How this is done is shown in Section 5.1, using a real data example.
2. *Dimensionality reduction* is sometimes done for performance reasons due to the high number of bands. The most common way is principal component analysis (PCA). Sometimes other, more theoretically suitable techniques are used to reduce dimensionality. For example, the minimum noise fraction (MNF) transform which finds uncorrelated linear transformations that maximize the signal-to-noise ratio. We do not consider these techniques in this thesis and always use the complete dataset.
3. *Unmixing/Inversion* is the last step in which we identify the endmembers in the scene and the fractional abundances in each pixel. This step requires that we understand how endmember spectra are mixed in each pixel. In this thesis we assume the linear mixing model (LMM, cf. Section 1.2.1). Unmixing is thus the reverse process of Eq. 1.1 under non-negativity and sum-to-one constraints. If we consider the matrix formulation as in Eq. 1.3, then this process is a form of non-negative matrix factorization (NMF) with the additional sum-to-one constraint for the abundances.

The unmixing step is at the core of this thesis. It is thus appropriate to start with an overview of the present unmixing methods.

1.3.1 Unmixing methods

Bioucas-Dias e.a. [8] give a recent overview of unmixing algorithms. They divide the approaches in three general classes: geometric, statistical and sparse. The following gives a brief overview.

Geometrical methods

The algorithms in this class are based on the minimal volume spanned by the endmembers. Broadly speaking, they differ in whether pure pixels are assumed or not. A popular algorithm of the former kind is the N-FINDR algorithm [26]. It is based on the idea that the volume defined by a simplex formed by the purest pixels, is larger than any other volume defined by any other combination of pixels. This algorithm finds the set of pixels defining the largest volume by inflating a simplex inside the data. It performs well under full mixing, but not so much under partial mixing (cf. Fig. 2.1a), which brings us to the more general class of geometrical volume constrained approaches. These approaches seek an endmember matrix \mathbf{E} that minimizes the volume of the simplex defined by its columns. Since the pure pixels are no longer assumed, this results in a more difficult nonconvex optimization problem. One such algorithm is the ICE algorithm [5], [6], [7] which is extensively discussed in Chapter 2.

Statistical methods

Statistical methods allow more precise modeling but come with a price of higher computational complexity. Especially the Bayesian approaches are popular. They have the ability to model statistical variability and to impose priors that can constrain solutions to physically meaningful ranges, and even regularize solutions. Common estimators are maximum a posteriori (MAP) estimators. Due to the complexity of resulting joint posterior, Markov chain Monte Carlo (MCMC) algorithms are used. These algorithms mainly differ by the choice of priors assigned to the unknown parameters.

In Chapter 3 we cover in great detail the BayesNMF-Vol [1], [3] algorithm which can be seen as a translation of the ICE algorithm into a Bayesian framework. Published in 2009, it is also one of the first Bayesian frameworks for hyperspectral unmixing. In the following years there was lot of research on this topic. For example this relatively new paper by Halimi e.a. [17] from 2015 is several layers more complex than the one described here. Halimi does not only incorporate the non-negativity and sum-to-one constraints into the model (which is an endeavor on its own), but goes further to model spatial information *and* spectral endmember variability (SEV) which is intrinsic [24] to remotely sensed spectral images.

Sparse regression methods

These methods are mostly semi-supervised in the sense that a large amount of spectra is fed into the model. The sparsity is related to the amount of endmembers in the solution (cf. [25]). There are also sparse unsupervised methods. We do not consider sparse methods in this thesis.

1.3.2 Spatial information

Most of the algorithms in the three classes ignore the spatial information that could improve the unmixing process. Bioucas-Dias et al. [8] note though that incorporating this information into a model has recently motivated the development of a new class of algorithms. These algorithms exploit both the spatial and spectral features contained in an image. We do not cover these algorithms, but we do extend the ICE algorithm with spatial regularization. To our knowledge this has not been done before.

1.4 Applications

Higher spectral resolution of hyperspectral cameras enables material identification through spectroscopic analysis. This facilitates countless applications that require identifying materials in scenarios unsuitable for classical spectroscopic analysis. For example, hyperspectral cameras such as AVIRIS contribute significantly to earth observation and remote sensing. In Chapter 5 we base our analysis on one such real dataset.

The applications are not limited to remote sensing. Hyperspectral images can be found in any other domain where material identification through spectroscopic analysis is desired. Bioucas-Dias [8] refers to several domains such as food safety, pharmaceutical process monitoring and quality control, biomedical, industrial, biometric, and forensic applications.

1.5 Notation and conventions

Variable scalars are written in lowercase italics (a, b), column vectors in lowercase boldface italics (\mathbf{a}, \mathbf{b}) and matrices in uppercase boldface letters (\mathbf{A}, \mathbf{B}). Transposition is denoted by the superscript T ($\mathbf{a}^T, \mathbf{A}^T$).

For notational brevity, the columns or rows of a $N \times M$ matrix \mathbf{A} are represented as a column vector \mathbf{a} where the subscript denotes the origin: $\mathbf{a}_n := (\mathbf{A}_{n,:})^T$ and $\mathbf{a}_m := \mathbf{A}_{:,m}$. Similarly, $a_{nm} := \mathbf{A}_{n,m}$ is the (n, m) -th component of matrix \mathbf{A} . We use tilde ($\tilde{\cdot}$) to denote removed indices such that $\mathbf{A}_{\tilde{a},\tilde{b}}$ denotes the \mathbf{A} matrix with row a and column b removed. Similarly $\mathbf{a}_{n\tilde{m}}$ denotes the vector \mathbf{a}_n extracted from \mathbf{A} with the m -th element removed.

The probability density function of a conventional distribution we denote in calligraphic style, like \mathcal{N} for normal distribution. For other distributions we use the generic lowercase letter f which is determined by its argument or subscript, or both. For example, $f_{X|Y}(x|y) = f(x|y) = f_{X|Y}$. The latter form we use mainly in text to refer to a specific function, while the middle form is used in equations where we keep the capital letters such as $f(X|Y)$. The former and most correct notation is omitted for the sake of brevity.

Some special matrices and vectors include the identity matrix \mathbf{I}_M with subscript denoting its order; a vector of ones $\mathbf{1}$, and the Moore-Penrose (MP) inverse of a matrix \mathbf{A} , which we write as \mathbf{A}^\dagger . The $\arg \min_{\mathbf{a}}$ of some function is denoted by \mathbf{a}^* .

Some important constants like the number of bands B , the number of endmembers M and the number of pixels N are written in uppercase italics. We use the same notation for important functions like the objective function L and the regularization functions V and S .

For matrix differential calculus we use the conventions laid out in [21]. Derivations are given only if they differ from the original papers.

1.6 Thesis organization and contributions

The following chapters are organized as follows. In Chapter 2 we explain Berman’s [5] ICE algorithm, which is perhaps the most typifying of the geometric class of algorithms. This results in a few small improvements: The closed form solution for \mathbf{E}^* in Section 2.3.4 and stabilization of the objective function in Section 2.4.1. Subsequently, in Section 2.5 we propose a spatial information extension for the ICE algorithm, which we dub ICE-S. The extra parameter ν of ICE-S allows for smoother abundance maps, but does not seem to have much effect on the resulting endmembers. We also notice multiple local minima in the objective function and try to map them by running ICE-S over its whole parameter range in Section 4.3.2.

In Chapter 3 we explain the Bayesian framework and the MCMC solution for hyperspectral unmixing, which is mainly due to Arngren e.a. [1]. Their work, back in 2009, was one of the pioneering in the field and their insights are fundamental for all subsequent Bayesian unmixing models. That is why we put a considerable amount of effort into understand it.

Subsequently, we run these three algorithms on a synthetic (Chapter 4) and a real (Chapter 5) dataset and compare the results. We end with a small conclusion in Chapter 6. In the Appendices one will find all the derivations that are omitted from the corpus, and are not included in the original papers. There, one will also find the link and explanation of the source code that accompanies this thesis.

Chapter 2

ICE algorithm and extensions

Berman’s iterated constrained endmembers (ICE) algorithm [5], [6], [7] is greatly inspired by Winters N-FINDR algorithm [26] which finds the M-simplex of *maximum* volume constrained to lie within the data cloud and Craig’s algorithm [13] which finds the *minimum* volume M-simplex enclosing the data cloud. The ICE algorithm finds a M-simplex which has Craig’s solution as its limiting case, and where the size of the simplex is controlled with a hyper-parameter μ . The comparison is shown in Fig. 2.1.

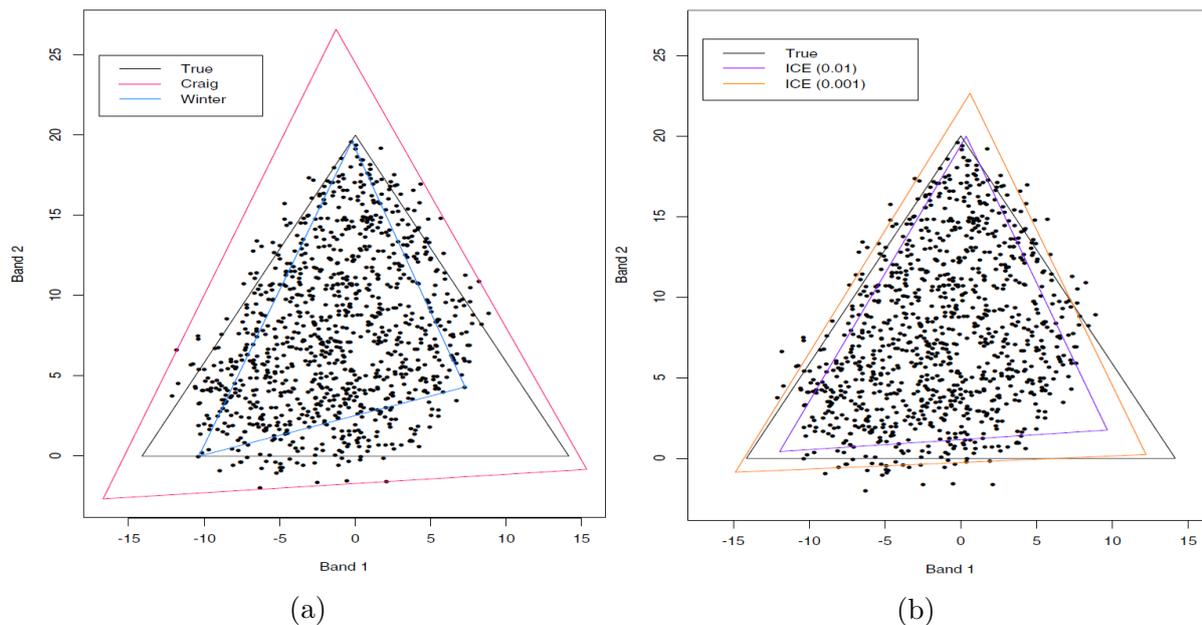


Figure 2.1: Comparison of Winter’s N-FINDR, Craig’s algorithm and Berman’s ICE algorithm for finding endmembers. Image taken from [7].

In Section 2.1 we explore the assumptions of the objective function central to the ICE algorithm. More precisely we derive it from LMM (Section 1.2.1) based on the Likelihood Principle. This is instructive as no such foundation has been found in the original papers. Once we learn that the objective function L is underspecified, we introduce in subsequent section a regularization term V which is proportional to the volume of the endmember simplex as a remedy. In Section 2.3 we explore in great detail every step in the so-called ALS algorithm which is used to minimize L . This Section also goes much deeper into implementation details than the original papers. Once done, we discuss two peculiarities of the ICE algorithm: omission of σ^2 and lurking imbalance due to large B – the number of

bands. And lastly in Section 2.5 we extend the ICE algorithm with spatial regularization and an extra hyper-parameter ν which controls it. To our knowledge this has not been tried before and is a creative contribution of this thesis.

2.1 Objective function

We start by rewriting \mathbf{x}_n from Eq. 1.1 in each component x_{nb} :

$$\begin{aligned} x_{nb} &= \mathbf{w}_n^T \mathbf{e}_b + \epsilon \\ &\sim \mathcal{N}(\mathbf{w}_n^T \mathbf{e}_b, \sigma^2) \end{aligned} \quad (2.1)$$

The probability density function of x_{nb} is thus:

$$f(x_{nb} | \mathbf{w}_n, \mathbf{e}_b, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2} \quad (2.2)$$

Assuming that all pixel components x_{nb} are independent we can write the joint distribution as follows:

$$\begin{aligned} f(x_{11}, \dots, x_{NB} | \mathbf{E}, \mathbf{W}, \sigma^2) &= f(x_{11} | \mathbf{w}_1, \mathbf{e}_1, \sigma^2) \times \dots \times f(x_{NB} | \mathbf{w}_N, \mathbf{e}_B, \sigma^2) \\ &= \prod_{n=1}^N \prod_{b=1}^B f(x_{nb} | \mathbf{w}_n, \mathbf{e}_b, \sigma^2) \\ &= \ell(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}) \end{aligned} \quad (2.3)$$

Now, from the Likelihood Principle all we have to do is maximize $\ell(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X})$. To facilitate the computation (i.e. to get rid of the product which is likely to cause numerical instability for large N and to get rid of the exponent), one usually uses the log-likelihood function $\mathcal{L} = \log(\ell)$. This is justified since the logarithm is a monotonically increasing function and thus preserves the maximum.

$$\begin{aligned} \mathcal{L}(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}) &= \log [\ell(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X})] \\ &= \sum_{n=1}^N \sum_{b=1}^B \log [f(x_{nb} | \mathbf{w}_n, \mathbf{e}_b, \sigma^2)] \\ &= -NB \log \left(\sqrt{2\pi\sigma^2} \right) - \frac{1}{2\sigma^2} \sum_{n=1}^N \sum_{b=1}^B (x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2 \end{aligned} \quad (2.4)$$

Now we are ready to define the objective loss function which we wish to minimize.

$$\begin{aligned} \arg \max_{\mathbf{E}, \mathbf{W}, \sigma^2} \mathcal{L}(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}) &= \arg \min_{\mathbf{E}, \mathbf{W}, \sigma^2} -\mathcal{L}(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}) \\ &= \arg \min_{\mathbf{E}, \mathbf{W}, \sigma^2} NB \log \left(\sqrt{2\pi\sigma^2} \right) + \frac{1}{2\sigma^2} \sum_{n=1}^N \sum_{b=1}^B (x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2 \end{aligned} \quad (2.5)$$

Note that Berman [5] seems to disregard σ^2 – without explanation¹ – and only considers \mathbf{E} and \mathbf{W} . In that case σ^2 is considered given and the additive constant $NB \log(\sqrt{2\pi\sigma^2})$ and positive factor $1/2\sigma^2$ can be removed since they have no impact on arg min.

$$\begin{aligned} \arg \max_{\mathbf{E}, \mathbf{W}} \mathcal{L}(\mathbf{E}, \mathbf{W} \mid \mathbf{X}) &= \arg \min_{\mathbf{E}, \mathbf{W}} \sum_{n=1}^N \sum_{b=1}^B (x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2 \\ &= \arg \min_{\mathbf{E}, \mathbf{W}} L(\mathbf{W}, \mathbf{E}, \sigma^2 \mid \mathbf{X}) \end{aligned} \quad (2.6)$$

The function $L(\mathbf{W}, \mathbf{E})$ we call the loss function or the objective function. At this point it is instructive to show various ways in which the objective function can be rewritten, which will prove useful in later manipulations.

$$\begin{aligned} L(\mathbf{E}, \mathbf{W} \mid \mathbf{X}) &= \sum_n^N \sum_b^B (x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2 = \sum_n^N \sum_b^B (x_{nb} - \hat{x}_{nb})^2 \\ &= \|\mathbf{X} - \mathbf{W}\mathbf{E}^T\|_F^2 = \text{tr} [(\mathbf{X} - \mathbf{W}\mathbf{E}^T)^T (\mathbf{X} - \mathbf{W}\mathbf{E}^T)] \\ &= \sum_n^N (\mathbf{x}_n - \mathbf{E}^T \mathbf{w}_n)^T (\mathbf{x}_n - \mathbf{E}^T \mathbf{w}_n) \\ &= \sum_b^B (\mathbf{x}_b - \mathbf{W}\mathbf{e}_b)^T (\mathbf{x}_b - \mathbf{W}\mathbf{e}_b) \end{aligned} \quad (2.7)$$

2.2 Regularized objective function

Note that objective function L in Eq. 2.7 is ill posed. It will result in the same minimum value as long as all the pixels are within the simplex formed by the vertices of the endmembers. In that case $x_{nb} - \hat{x}_{nb} = 0 \forall n, b$ and thus $L = 0$. That is why Berman [5] proposes a penalized version:

$$L_{reg}(\mathbf{E}, \mathbf{W}) = \frac{(1 - \mu)}{N} L + \mu V \quad (2.8)$$

Here V acts as a regularization term and is proportional to the volume of the end-member simplex. μ acts as a hyper-parameter of the model. Division by N is to make sure that L is independent of the amount of pixels used. To make V independent of the dimension of the simplex, the summed variance per band b of all the endmembers is used instead.

¹Actually, that is a bit too quick. Berman uses the ICE algorithm mainly on data which is transformed by means of the minimum noise fraction (MNF) transform. This makes the errors uncorrelated between MNF bands and with variance 1. The data can still be spatially correlated. Furthermore, MNF transform tends to produce errors that look Gaussian for low dimensions. [6]

$$\begin{aligned}
V &= \sum_{b=1}^B \widehat{Var}(e_b) \\
&= \sum_{b=1}^B \frac{\sum_{m=1}^M \left[e_{bm} - \frac{\sum_{m=1}^M e_{bm}}{M} \right]^2}{M-1} \\
&= \sum_{b=1}^B \frac{e_b^T (\mathbf{I}_M - \frac{\mathbf{1}\mathbf{1}^T}{M}) e_b}{M-1} \\
&= \frac{SSD}{M(M-1)}
\end{aligned} \tag{2.9}$$

Note that Eq. 2.9 gives a direct relation between V and the sum of squared distances (SSD) of the endmember simplex. A worked out derivation of this relationship can be found in the Appendix Eq. A.2. Note further that the SSD is also proportional to the size of the simplex. In summary, the advantage of V to the volume is threefold: it is independent of M , proportional to the volume and cheaper to compute.

Since L_{reg} is approximately independent of the sample size N and the number of endmembers M , it allows for just one value of μ for all datasets. Berman [5], [7] reports that a value of 0.01 to 0.05 usually gives reasonable solutions after analyzing several dozen real-world datasets.

2.3 Minimization

Minimization of L_{reg} is not straightforward. The dimensionality is extremely large. Suppose for example one has an image of $N=1000$ pixels with $B=100$ bands and one seeks $M=5$ endmembers. Since our function L_{reg} has matrices $\mathbf{E} \in \mathbb{R}_+^{B \times M}$ and $\mathbf{W} \in \mathbb{R}_+^{N \times M}$ as arguments, we would have a total of $B \times M + N \times M = 100500$ parameters to find! And even if we had a closed-form solution, we still would need to make sure the non-negativity and sum-to-one constraints are met.

The approach taken by Berman is to split the function in smaller subproblems and solve them one by one. Berman only gives the closed form solution for e_b^* and mentions that w_n^* can be found using quadratic programming as explained in [22, Chapter 16]. In the following sections we delve into details of how this is actually done.

2.3.1 Coordinate descent

Coordinate descent [27] is based on the idea that the minimization of a multi-variable objective function $f(\mathbf{x})$ can be achieved by minimizing it along one component of \mathbf{x} with respect to the remaining components at each iteration step. One thus reduces the difficult full minimization problem to a cycle of much simpler single-variable problems. One starts with a random initial value:

$$\mathbf{x}^0 = (x_1^0, \dots, x_n^0)$$

Then one iteratively solves the single component optimization problems

$$x_i^{k+1} = \arg \min_{y \in \mathbb{R}} f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, y, x_{i+1}^k, \dots, x_n^k)$$

where the superscript $k + 1$ denotes the current iteration cycle. At each step one does a line search for which $f(\mathbf{x}^k) \geq f(\mathbf{x}^{k+1})$ always holds. Overall we get

$$f(\mathbf{x}^0) \geq f(\mathbf{x}^1) \geq f(\mathbf{x}^2) \geq \dots$$

No significant improvement between two iteration cycles implies a stationary point is reached.

2.3.2 Block coordinate descent

An interesting feature of the objective function $L_{reg}(\mathbf{W}, \mathbf{E})$ in Eq. 2.8 is that if we only look at $L_{reg}(\mathbf{w}_n)$ and $L_{reg}(\mathbf{e}_b)$ with other parameters fixed, then the function in these parameters is convex. Thus we can minimize efficiently these blocks of parameters one at a time. Hence the name *block* coordinate descent. Now let us take a closer look at these functions and their form.

2.3.3 Minimize $L_{reg}(\mathbf{w}_n)$ with fixed \mathbf{E} and $\mathbf{W}_{\tilde{n}}$:

Assuming fixed \mathbf{E} we can rewrite L_{reg} from Eq. 2.8 in function of \mathbf{w}_n as follows:

$$\begin{aligned} L_{reg}(\mathbf{W}) &= \sum_{n=1}^N \left[\frac{(1-\mu)}{N} \|\mathbf{x}_n - \mathbf{E}\mathbf{w}_n\|_2^2 + \frac{\mu V}{N} \right] \\ &= \sum_{n=1}^N L_{n,reg}(\mathbf{w}_n) \end{aligned} \quad (2.10)$$

This shows that each $L_{n,reg}(\mathbf{w}_n)$ can be minimized separately since \mathbf{w}_n are independent of each other. Now we rewrite $L_{n,reg}(\mathbf{w}_n)$ to quadratic standard form:

$$\begin{aligned} L_{n,reg}(\mathbf{w}_n) &= \frac{(1-\mu)}{N} (\mathbf{x}_n - \mathbf{E}\mathbf{w}_n)^T (\mathbf{x}_n - \mathbf{E}\mathbf{w}_n) + \frac{\mu V}{N} \\ &= \left(\frac{(1-\mu)\mathbf{x}_n^T \mathbf{x}_n}{N} + \frac{\mu V}{N} \right) - \mathbf{w}_n^T \frac{2(1-\mu)\mathbf{E}^T \mathbf{x}_n}{N} \\ &\quad + \frac{1}{2} \mathbf{w}_n^T \frac{2(1-\mu)\mathbf{E}^T \mathbf{E}}{N} \mathbf{w}_n \\ &= C^{te} + \mathbf{w}_n^T \mathbf{a} + \frac{1}{2} \mathbf{w}_n^T \mathbf{H} \mathbf{w}_n \end{aligned} \quad (2.11)$$

Note that the constant term plays no role in minimization. Thus we need to solve the following constrained quadratic problem:

$$\begin{aligned} &\arg \min_{\mathbf{w}_n} \mathbf{w}_n^T \mathbf{a} + \frac{1}{2} \mathbf{w}_n^T \mathbf{H} \mathbf{w}_n \\ &\text{subject to } \sum_{m=1}^M w_{nm} = 1 \text{ and } w_{nm} \geq 0, \quad m = 1, \dots, M. \end{aligned} \quad (2.12)$$

For numerical stability (in case N is large) we can remove the factor $2(1-\mu)/N$ from both \mathbf{a} and \mathbf{H} without influencing the minimum.

Minimizing $L_{reg}(\mathbf{w}_n)$ under constraints

Nocedal [22, Ch. 16] describes two approaches for solving such a problem. The first is the *active set method* which is inspired by the simplex method from linear programming. In the simplex method, a solution is searched for along the vertices of the feasible polytope. Each vertex of such a polytope is a point where all equality, and a subset of inequality constraints, are active. An inequality constraint such as $w_{nm} \geq 0$ is considered active if it strictly equals zero at that point.

Quadratic problems are more complicated since the solution is not necessarily one of the vertices, nor does it necessarily lie on the boundary of the polytope. Nonetheless, searching along the boundary of the polytope is the essence of the active set algorithms. For this particular problem we have implemented the *primal* active set method for convex quadratic problem as described in [22, p. 472].

A second approach is the so called *interior-point method*. This method is polynomial time and approaches the solution through the interior of the feasible polytope rather than working its way around the boundary as the active set method does.

Convexity of $L_{reg}(\mathbf{w}_n)$

$L_{n,reg}(\mathbf{w}_n)$ is convex and our minimum is global if \mathbf{H} is positive semidefinite. The minimum is unique if \mathbf{H} is positive definite. We know that $\mathbf{x}^T \mathbf{E}^T \mathbf{E} \mathbf{x} = (\mathbf{E} \mathbf{x})^T \mathbf{E} \mathbf{x} = \|\mathbf{E} \mathbf{x}\|_2^2 \geq 0$ and thus that \mathbf{H} is positive semidefinite. If $\ker(\mathbf{E}) = \{\mathbf{x} \mid \mathbf{E} \mathbf{x} = \mathbf{0}\} = \{\mathbf{0}\}$, then the columns of \mathbf{E} are linearly independent and strict inequality holds for $\mathbf{x} \neq \mathbf{0}$. Thus \mathbf{H} is positive definite only if the endmembers (i.e. columns) of \mathbf{E} are linearly independent! Since $\mathbf{E} \in \mathbb{R}_+^{B \times M}$, the latter can *only* be if $M \leq B$. So by having $M \ll B$ we increase the chance of having an unique solution for \mathbf{w}_n . The case where $M > B$ is common in synthetic images where we limit $B = 2$ and $M = 3$ for representational purposes as in Fig. 1.4. In such cases \mathbf{H} is positive semidefinite. If this were an unconstrained problem, then we would be sure that multiple solutions for \mathbf{w}_n exist. In the constrained case this is more complicated. The set of minima for semi-definite \mathbf{H} might as well lie completely outside of the feasible set, so we might still have an unique solution under active constraints! In other words, in the constrained case, the requirements for a unique solution are more loose. Specifically [22, Lemma 16.1, p.452], an active set of constraints \mathcal{A} with its corresponding matrix \mathbf{A} has a unique solution if $\mathbf{Z}^T \mathbf{H} \mathbf{Z}$ is positive definite where $\ker(\mathbf{A}) = \mathbf{Z}$. So even if \mathbf{H} is positive semi-definite, we still have a chance that the constrained solution is unique.

Implementation peculiarities

Note that we *could* run the active set algorithm on the full matrix \mathbf{W} since $L_{reg}(\mathbf{W})$ is quadratic and convex, but this is dangerous. Nocedal [22, p. 388] notes that the complexity of an active set method is exponential, i.e. in the linear case there exist pathological cases for which the simplex method visits every single vertex before reaching the optimal point. Therefore it is safer to limit the dimensionality of the problem to just $L_{reg}(\mathbf{w}_n)$.

Another implementation peculiarity is the case when $\mathbf{Z}^T \mathbf{H} \mathbf{Z}$ is positive semi-definite. In that case the Karush-Kuhn-Tucker (KKT) matrix is singular and thus there are multiple solutions for \mathbf{w}_n . A general approach for solving such singular case is to use the Moore-Penrose (MP) pseudo-inverse which gives a solution \mathbf{x}_{min} having the smallest norm

of all possible solutions \mathbf{x} . An important thing to note here is that $\mathbf{x}_{min} = (\mathbf{w}_n, \boldsymbol{\lambda})$, where $\boldsymbol{\lambda}$ is the vector of Lagrange multipliers. So although \mathbf{x}_{min} has the minimal norm, that does not imply that \mathbf{w}_n also has the minimal norm of all possible solutions.

Another issue with MP inverse is its calculation. It is usually computed using singular value decomposition (SVD). Thus $\mathbf{A}^\dagger = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^T$ where \mathbf{D}^{-1} is the diagonal matrix of reciprocal eigenvalues. Very small eigenvalues cause MP inverse instability. On the other hand, removing small but non-zero eigenvalues causes inexact MP inverse which eventually results in solutions that might not exactly match the constraints. So although rare, these situations occur, especially when $M > B$ such as in synthetic images, and might cause errors of order 1e-8 or even higher.

During the process of implementation and testing we have found an inconsistency in the active set algorithm for convex QP as described by Nocedal [22, p.472]. Where it says “obtain W_{k+1} by adding one of the blocking constraints to W_k ” it should be “obtain W_{k+1} by adding *the* blocking constraint which corresponds with the smallest α_k to W_k ”. This is essential as picking a random blocking constraint when there are multiple blocking constraints will, result in a wrong minimum for a well-chosen initial feasible point x_0 .

2.3.4 Minimize $L_{reg}(\mathbf{e}_b)$ with fixed \mathbf{W} and $\mathbf{E}_{\tilde{b}}$:

Assuming fixed \mathbf{W} , using $\|\mathbf{X} - \mathbf{W}\mathbf{E}^T\|_F^2 = \sum_{b=1}^B \|\mathbf{x}_b - \mathbf{W}\mathbf{e}_b\|_2^2$ and substituting the result for V from Eq. 2.9 into Eq. 2.8 we get:

$$\begin{aligned} L_{reg}(\mathbf{E}) &= \sum_{b=1}^B \left[\frac{(1-\mu)}{N} \|\mathbf{x}_b - \mathbf{W}\mathbf{e}_b\|_2^2 + \mu \frac{\mathbf{e}_b^T (\mathbf{M}\mathbf{I}_M - \mathbf{1}\mathbf{1}^T) \mathbf{e}_b}{M(M-1)} \right] \\ &= \sum_{b=1}^B L_{b,reg}(\mathbf{e}_b) \end{aligned} \quad (2.13)$$

This shows that each $L_{b,reg}(\mathbf{e}_b)$ can be minimized separately since \mathbf{e}_b are independent of each other. Now we rewrite $L_{b,reg}(\mathbf{e}_b)$ in quadratic standard form:

$$\begin{aligned} L_{b,reg}(\mathbf{e}_b) &= \frac{(1-\mu)\mathbf{x}_b^T \mathbf{x}_b}{N} - \mathbf{e}_b^T \frac{2(1-\mu)\mathbf{W}^T \mathbf{x}_b}{N} \\ &\quad + \frac{1}{2} \mathbf{e}_b^T \left(\frac{2(1-\mu)\mathbf{W}^T \mathbf{W}}{N} + \frac{2\mu(\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M)}{(M-1)} \right) \mathbf{e}_b \\ &= C^{te} + \mathbf{e}_b^T \mathbf{c} + \frac{1}{2} \mathbf{e}_b^T \mathbf{G} \mathbf{e}_b \end{aligned} \quad (2.14)$$

Convexity of $L_{reg}(\mathbf{e}_b)$

The Hessian matrix \mathbf{G} is symmetric since it is a sum of two symmetric matrices. \mathbf{G} is also positive semi-definite. To see this we have to look at its matrix parts. $\mathbf{e}^T \mathbf{W}^T \mathbf{W} \mathbf{e} = \|\mathbf{W}\mathbf{e}\|_2^2 \geq 0$ so $\mathbf{W}^T \mathbf{W}$ is positive semi-definite. That $(\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M)$ is also positive semi-definite follows from $\mathbf{x}^T (\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M) \mathbf{x} = \mathbf{x}^T (\mathbf{I}_M - \mathbf{v}\mathbf{v}^T) \mathbf{x} = \|\mathbf{x}\|_2^2 - \|\mathbf{v}^T \mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2 \|\mathbf{v}\|_2^2 - \|\mathbf{v}^T \mathbf{x}\|_2^2 \geq 0$ where $\mathbf{v} = \mathbf{1}/\sqrt{M}$ and where the last inequality is due to the Cauchy-Schwarz inequality. The consequence is that $L_{b,reg}(\mathbf{e}_b)$ is convex. In practice, \mathbf{G} is positive definite because to be semi-definite the columns of \mathbf{W} have to be dependent and since $N \gg M$ this is highly unlikely. At the same time the $\ker(\mathbf{W}^T \mathbf{W})$ has to be the

same as $\ker(\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M) = \{\mathbf{x} \mid x_1 = x_2 = \dots = x_M\}$ which is even more unlikely. So in the calculation we can safely use \mathbf{G}^{-1} instead of \mathbf{G}^\dagger .

Closed solution for \mathbf{e}_b

The convexity of the problem thus assures us that solving for \mathbf{e}_b in equation $\partial L_{b,reg}(\mathbf{e}_b)/\partial \mathbf{e}_b^T = \mathbf{c}^T + \mathbf{e}_b^T \mathbf{G} = 0$ will result in the global minimum:

$$\arg \min_{\mathbf{e}_b} L_{b,reg}(\mathbf{e}_b) = \left[\mathbf{W}^T \mathbf{W} + \lambda \left(\mathbf{I}_M - \frac{\mathbf{1}\mathbf{1}^T}{M} \right) \right]^{-1} \mathbf{W}^T \mathbf{x}_b \quad (2.15)$$

with $\lambda = N\mu/[(1-\mu)(M-1)]$.

But are we not a bit too fast? What about the constraint $e_{bm} \geq 0, \forall b, m$? It is a peculiarity in Berman's ICE algorithm that it is not enforced. But while it is not enforced, it is very likely to be satisfied since the endmembers will be within the pixel cloud. And since the pixel components are all greater or equal to zero, it is safe to assume that the endmember components will also be.

Closed solution for \mathbf{E}

Although Berman [5], [6], [7] does not mention it explicitly, $\arg \min_{\mathbf{E}} L_{reg}(\mathbf{E})$ can easily be solved in one go by replacing \mathbf{x}_b with \mathbf{X} and taking the transpose of the result in Eq. 2.15, resulting in:

$$\arg \min_{\mathbf{E}} L_{reg}(\mathbf{E}) = \mathbf{X}^T \mathbf{W} \left[\mathbf{W}^T \mathbf{W} + \lambda \left(\mathbf{I}_M - \frac{\mathbf{1}\mathbf{1}^T}{M} \right) \right]^{-1} \quad (2.16)$$

This result follows directly from Eq. 2.15 once we recognize that the $L_{b,reg}(\mathbf{e}_b)$ are independent. A more formal derivation based on matrix calculus laws from [21] is given in Appendix A.2.

2.3.5 The ICE algorithm

In the previous two sections we computed the parameters \mathbf{E}^{k+1} and \mathbf{W}^{k+1} for one whole cycle of a coordinate descent. We thus have that $L_{reg}(\mathbf{E}^k, \mathbf{W}^k) \geq L_{reg}(\mathbf{E}^{k+1}, \mathbf{W}^{k+1})$. Since at each step we are minimizing a quadratic problem, this approach is also called alternating least squares (ALS). We now describe the full algorithm.

The iteration process is stopped once successive L_{reg} values are close enough. Berman [6] reports that a ratio larger than 0.99999 gives sufficiently stable endmember estimates. Now the obvious question is: does this converge to a global minimum?

2.3.6 Non-Convexity of L_{reg}

The global minimum is achieved by coordinate descent only if the objective function is convex. Unfortunately, L_{reg} as defined in Eq. 2.8 is not convex [6]. There is no formal proof for this statement. But one can demonstrate this by using different starting points in the ICE algorithm. For example, if the starting endmember simplex is outside of the data cloud, the solution will sometimes tend to equal endmembers in the center of the data cloud. This seems to be a local minimum. An other local minimum seems to be when the endmembers lie on one straight line within the data cloud. From our conducted

Algorithm 1 (ICE algorithm)

```
 $k \leftarrow 1$   
 $L^{(k)} \leftarrow$  maximum machine number  
Generate a random matrix  $\mathbf{W}^{(k)}$   
repeat  
   $k \leftarrow k + 1$   
   $\mathbf{E}^{(k)} \leftarrow \arg \min_{\mathbf{E}} L_{reg}(\mathbf{E} \mid \mathbf{W}^{(k-1)})$  as in Eq. 2.16  
  for  $n = 1$  to  $N$  do  
     $\mathbf{w} \leftarrow \arg \min_{\mathbf{w}_n} L_{n,reg}(\mathbf{w}_n \mid \mathbf{E}^{(k)})$  as in Eq. 2.12  
     $\mathbf{W}_n^{(k)} \leftarrow (\mathbf{w})^T$   
  end for  
   $L^{(k)} \leftarrow L_{reg}(\mathbf{W}^{(k)}, \mathbf{E}^{(k)})$  as in Eq. 2.8  
until  $L^{(k)}/L^{(k-1)} \geq tol$   
return  $\mathbf{E}^{(k)}$  and  $\mathbf{W}^{(k)}$  as solution;
```

experiments on synthetic data, it also seems better to start the algorithm with a random \mathbf{W} – instead of a random \mathbf{E} – to end up at the global minimum.

Berman’s experience is that endmembers which are pure or almost pure are consistently found, almost independently of the starting point, in the sense that there is little variation in the solutions. Less pure endmembers are found with greater variation.²

Sampling the solutions while initializing the algorithm with random \mathbf{W} also shows that the chance to end up in a local minima increases with decreasing μ . We will return to this point in Chapter 4 (Fig. 4.2a), when we sample the solutions over the whole parameter space of μ .

2.4 Model peculiarities

2.4.1 L dependence on B

During our testing the $\arg \min_{\mathbf{w}_n} L_{n,reg}(\mathbf{w}_n)$ proved instable for bands $B \geq 100$. This was due to $\mathbf{E}^T \mathbf{E}$ in effect having very large coefficients in case \mathbf{X} was not normalized. This in turn resulted in the KKT matrix having very large and very small eigenvalues (due to small coefficients of active constraint matrix \mathbf{A}), making the KKT matrix computationally singular. This results in numerical instability in the solution of \mathbf{w}_n^* with the effect of not complying to the sum-to-one constraint. That is because the small eigenvalue, corresponding to the equality constraint, gets dropped in the MP inverse calculation because it is too small compared to the others.

The instability can be traced back to L and V being both independent of M and N , but not of B . So the larger the B , the larger L_{reg} will tend to be. This can be mitigated by dividing L_{reg} in Eq. 2.8 with B . Our new improved objective function thus becomes:

$$L_{imp,B}(\mathbf{W}, \mathbf{E}) = \frac{(1 - \mu)}{NB} L + \frac{\mu V}{B} \quad (2.17)$$

The effect is that this divides the factors \mathbf{a} and \mathbf{H} in Eq. 2.12 by B mitigating the large eigenvalues without affecting the solution.

²cf. comments at: https://www.researchgate.net/publication/229884791_ICE_A_new_method_for_the_multivariate_curve_resolution_of_hyperspectral_images

$$\mathbf{a} = \frac{\mathbf{E}^T \mathbf{x}_n}{B}, \quad \mathbf{H} = \frac{\mathbf{E}^T \mathbf{E}}{B} \quad (2.18)$$

Note that this adjustment has no effect on \mathbf{E}^* .

2.4.2 Disregard of σ^2

We already noted that Berman does not consider σ^2 in the likelihood Eq. 2.4. For the sake of example, let us add it into L_{reg} and see what happens.

$$\begin{aligned} L_{imp,\sigma^2}(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}) &= -\frac{(1-\mu)}{NB} \mathcal{L}(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}) + \frac{\mu V}{B} \\ &= \frac{(1-\mu)}{2\sigma^2 NB} \|\mathbf{X} - \mathbf{W}\mathbf{E}^T\|_F^2 + (1-\mu) \log(\sqrt{2\pi\sigma^2}) + \frac{\mu V}{B} \end{aligned} \quad (2.19)$$

To minimize this we have to make slight adjustments to the ICE algorithm. \mathbf{W}^* remains the same. For computing \mathbf{E}^* we use the Eq. 2.16 with adjusted $\lambda = 2\sigma^2 N\mu / [(1-\mu)(M-1)]$. And finally, $\sigma^{2*} = \|\mathbf{X} - \mathbf{W}\mathbf{E}^T\|_F^2 / (NB)$.

Although all the subproblems are convex, the function L_{reg} is definitely not. From practical experiments, it seems that there are many local minima.

2.5 Extension with spatial regulation (ICE-S)

2.5.1 Motivation

Suppose you were given two abundance maps of some mineral as in Fig. 2.2. Both abundance maps contain the same data, but the abundances are ordered differently. If these two images were both the output for the same material, one would tend to think that in the right image the algorithm was somehow malconfigured. The left image shows some configuration, order or a hidden manifold. In the right image the abundances are randomly spread throughout the image which makes it less likely.

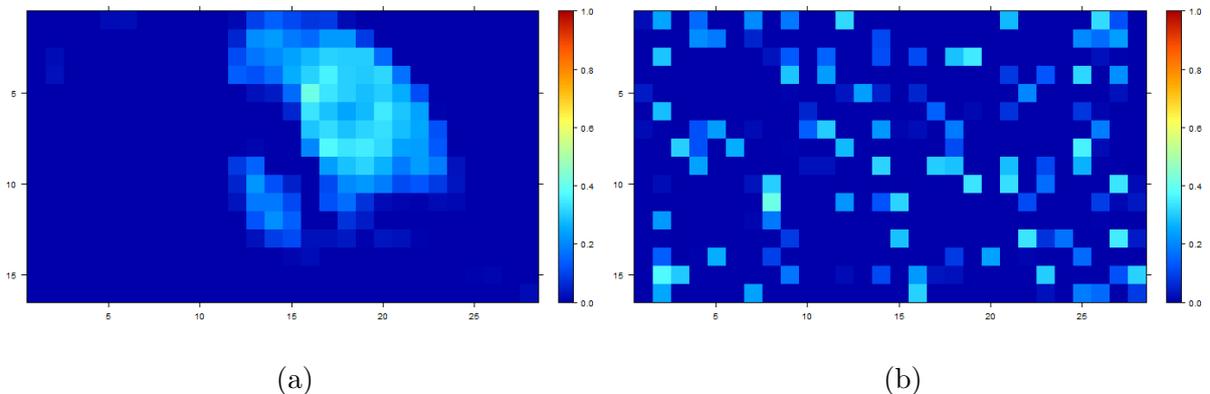


Figure 2.2: Spatial abundance of a mineral Alunite. The left image (a) is a copy of Fig. 5.7a – the original abundance. The right image (b) has the same abundance data randomly ordered.

We want a regularization that favors the left image, and disfavors the right image. There are various ways one could do this. The approach chosen here is to model each pixel with its four neighboring pixels as cliques and compute their variance.

2.5.2 New regularization function S

Thus the variance of an abundance w_{nm} and its neighbors on band m we define as:

$$S_{nm} = \text{Var}(\mathbf{U}_{:m}^{[n]}) \text{ where } \mathbf{U}^{[n]} = \begin{bmatrix} \mathbf{w}_{n-H}^T \\ \mathbf{w}_{n-1}^T \\ \mathbf{w}_n^T \\ \mathbf{w}_{n+1}^T \\ \mathbf{w}_{n+H}^T \end{bmatrix} \quad (2.20)$$

where H denotes the height of the image and is important for the calculation of adjacent pixels. Thus the S_{nm} measures the variance between abundances $\{w_{n-H,m}, w_{n-1,m}, w_{n,m}, w_{n+1,m}, w_{n+H,m}\}$ as shown in Fig. 2.3. The vector $\mathbf{U}_{:m}^{[n]}$ contains all the necessary information around abundance w_{nm} for this calculation. The total variance is simply the sum of all S_{nm} . Thus $S = \sum_{n=1}^N \sum_{m=1}^M S_{nm}$.

$w_{0,m}$				
\vdots		$w_{n-1,m}$		
	$w_{n-H,m}$	$w_{n,m}$	$w_{n+H,m}$	
		$w_{n+1,m}$		\vdots
				$w_{N,m}$

Figure 2.3: A 5×5 grid composed of endmember abundances for band m . The central 5 abundances are the components of $\mathbf{U}_{:m}^{[n]}$ defined in Eq. 2.20.

2.5.3 New objective function L_{VS}

The proposed new regularized objective function:

$$L_{VS}(\mathbf{W}, \mathbf{E} \mid \mathbf{X}) = \frac{(1-\mu)}{NB} \|\mathbf{X} - \mathbf{W}\mathbf{E}^T\|_F^2 + \mu \left[\frac{\nu}{B} V + \frac{(1-\nu)}{NM} S \right] \quad (2.21)$$

where ν acts as a regularization parameter that distributes the regularization between volume V and spatial information S . Note that we also divide S by N and M , which is to make it of same magnitude as V .

2.5.4 Minimization

The fact that we use variance ensures the quadratic form of $S_n(\mathbf{W})$. This fits nicely into the existing ICE algorithm.

Derivation of $L_{VS}(\mathbf{W})$

We start from Eq. 2.21. Since we are using coordinate descent, we focus on solving \mathbf{w}_n one at a time, just as we did previously. This time we have S which we need to write as a function of \mathbf{w}_n . The exact derivation is given in Appendix A.3.

$$S_n(\mathbf{w}_n) = \sum_{k \in \mathcal{K}_n} \frac{\mathbf{w}_n^T \mathbf{w}_n}{K_k} - 2\mathbf{w}_n^T \sum_{k \in \mathcal{K}_n} \frac{(\mathbf{U}_{\tilde{n}:}^k)^T \mathbf{1}}{K_k(K_k - 1)} + C^{te} \quad (2.22)$$

The set $\mathcal{K}_n = \{n - H, n - 1, n, n + 1, n + H\}$ is defined here as the set of all adjacent numbers of abundance-pixel n , including n itself. Note that cardinality $K_k = |\mathcal{K}_k|$ is a variable and not always equal to 5. Using this result and Eq. 2.21 we can derive $L_{n,VS}$ for use in minimization:

$$\begin{aligned} L_{n,VS}(\mathbf{w}_n) &= \frac{(1 - \mu)}{NB} \|\mathbf{x}_n - \mathbf{E}\mathbf{w}_n\|_2^2 + \frac{\mu}{N} \left(\frac{\nu}{B} V + \frac{(1 - \nu)}{M} S_n \right) \\ &= C^{te} - \mathbf{w}_n^T \left[2 \frac{(1 - \mu) \mathbf{E}^T \mathbf{x}_n}{NB} + 2 \frac{\mu(1 - \nu)}{NM} \sum_{k \in \mathcal{K}} \frac{(\mathbf{U}_{\tilde{n}:}^k)^T \mathbf{1}}{K_k(K_k - 1)} \right] \\ &\quad + \frac{1}{2} \mathbf{w}_n^T \left[2 \frac{(1 - \mu) \mathbf{E}^T \mathbf{E}}{NB} + 2 \frac{\mu(1 - \nu)}{NM} \sum_{k \in \mathcal{K}} \frac{\mathbf{I}_M}{K_k} \right] \mathbf{w}_n \\ &= C^{te} + \mathbf{w}_n^T \mathbf{a}_S + \frac{1}{2} \mathbf{w}_n^T \mathbf{H}_S \mathbf{w}_n \end{aligned} \quad (2.23)$$

The factor $2/N$ can be removed since it does not affect the outcome. Now that we have \mathbf{a}_S and \mathbf{H}_S we can use the same minimization technique as for Eq. 2.12 to minimize it.

2.5.5 Derivation of $L_{VS}(\mathbf{E})$

There are hardly any changes here compared to Eq. 2.15. S depends on \mathbf{W} so it is a constant and has no effect on \mathbf{E}^* . The only thing that gets added to Eq. 2.15 is ν such that $\lambda = N\mu\nu/[(1 - \mu)(M - 1)]$.

2.6 Evaluation

Fig. 2.4 shows how spatial regulation influences the abundance maps. The more spatial regulation we apply through the $(1 - \nu)$ hyper-parameter – while holding $\mu\nu$ (i.e. volume regularization) constant – the more smooth the abundances. We do not see any improvement in endmember estimation, as Fig. 2.4d shows.

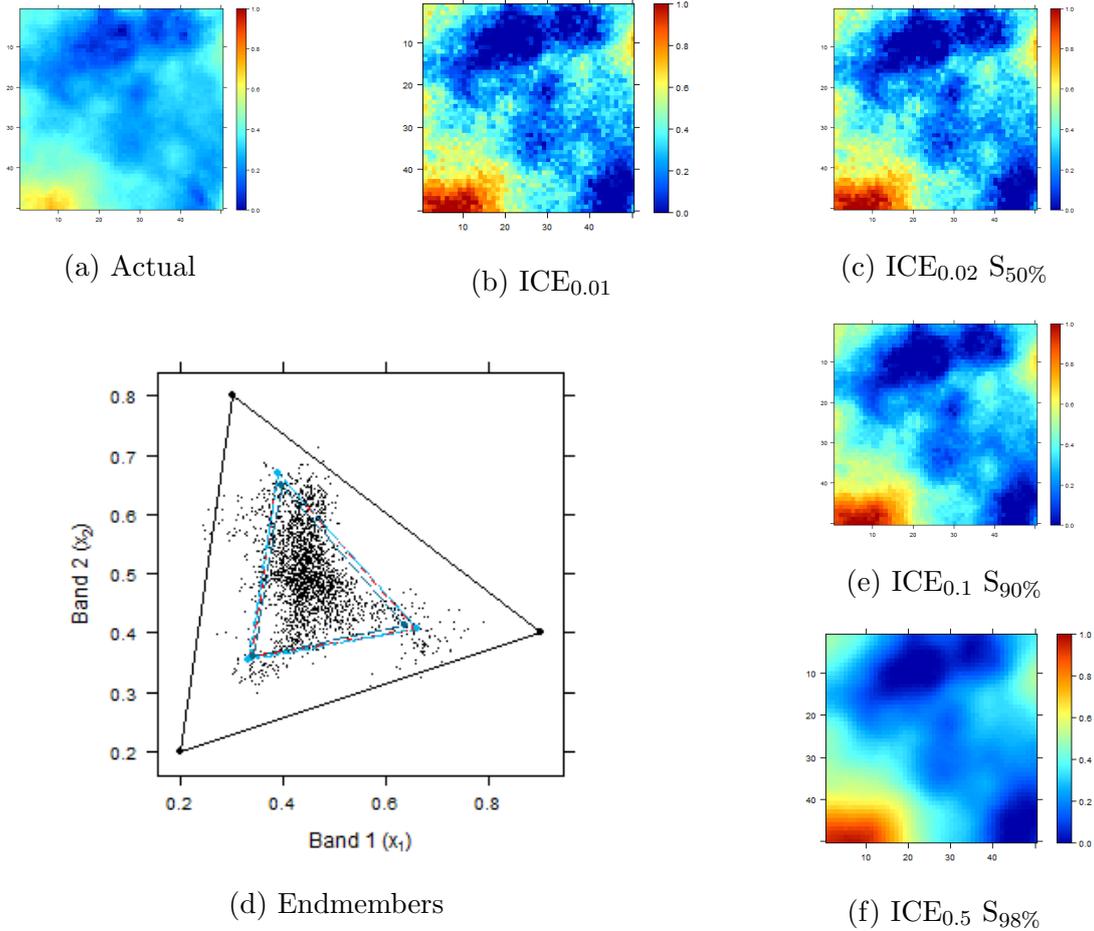


Figure 2.4: These images are taken from analysis in Chapter 4, Fig 4.3. Fig. (a) is the real abundance map. Fig. (b) is estimated with the ICE algorithm using $\mu = 0.01$. Note that 0.01 is suggested by Berman [7] as the optimal parameter based on analysis of various real datasets. Fig. (c) is estimated with the spatially extended ICE algorithm using $\mu = 0.02$ and $\nu = 0.5$ which results in weight 0.01 for V and 0.50 for S . Fig. (e-f) are similar, where weight for V is held constant with $\mu\nu = 0.01$. These images are merely to illustrate the effect of spatial regularization on the abundance maps. Fig. (d) compares the resulting endmembers (red for $\text{ICE}_{0.01}$, light to dark blue for the amount of regularization $S_{50\%}$, $S_{90\%}$, $S_{98\%}$ respectively while holding $\mu\nu = 0.01$) with the true endmembers (black). The endmember corresponding with the abundance maps is the bottom left one.

2.7 Conclusion

In this chapter we have seen how the objective function of the ICE algorithm can be derived from the LMM model using the Likelihood Principle. Subsequently we have thoroughly analyzed the minimization process of the objective function using coordinate descent. We made one contribution in this regard by dividing the objective function with B , thus improving its stability – especially for cases with unnormalized data. We also calculated the closed form solution for \mathbf{E}^* . Subsequently, driven by the idea that there is some order in abundances maps of materials, we have proposed and implemented a spatial regularization extension to the ICE algorithm, called ICE-S. This algorithm has one extra hyper-parameter $1 - \nu$ which controls the spatial regularization. The higher

we set this parameter, the smoother the abundance maps. The results also show that this parameter has no significant effect on endmember estimates – compared to the ICE algorithm under the same volume regularization weight.

Chapter 3

BayesNMF-Vol algorithm

A common approach to volume-constrained hyperspectral unmixing is to build a regularized objective function (like L_{reg} in Eq. 2.8) and solve for \mathbf{E}^* by numerical optimization. Another approach is to build a probability model and treat the extraction of endmembers \mathbf{E} as a Bayesian inference problem. This approach is based on [1], [3] and requires the definition of a likelihood, priors and most of all, a sound Belief (or Bayesian) Network.

Our primary focus here, like in the previous chapter, is understanding the principles and limitations of the approach. That means that we go further into details than the original papers. One essential point that is lacking in the original papers is the definition of a sound Belief Network. This step is essential because it defines the form of the joint distribution which in turn allows us to simplify calculations with Bayes rule. That is why we start this chapter with the definition of a Belief Network that underpins the BayesNMF-Vol model.

Then, in Section 3.2 we define the likelihood and priors in accordance with the original papers. We are mainly interested in the maximum a posteriori (MAP) estimators for this model. One reason is because there is a relationship between the MAP estimators and the maximum likelihood estimators (MLE) as used in the ICE Algorithm, which in turn allow direct comparison. We touch briefly on this relationship in Section 3.3.

The need for a posterior distribution calls for more complex MCMC techniques such as a Gibbs sampler on which we elaborate in Section 3.4. The original papers are accompanied with a technical paper [2] in which all the needed conditional distributions are derived. We refer to those when needed. In other cases where our own derivations differ (slightly), we put them in the Appendix. This is most notably the case for the \mathbf{W} -sampler in Section 3.4.4 where we are confronted with a degenerate distribution. In Section 3.2.4 we also look into the derivation and implementation of Jeffreys' prior which is a limiting case of the Inverse-Gamma distribution.

Lastly, in Section 3.5 we investigate some model peculiarities. The most notable and somewhat unexpected one is the intrinsic volume regularization of this Bayesian approach, due to the choice of the \mathbf{W} -prior. Several authors [16], [3] have a mathematical explanation for this phenomenon, while we give a more conceptual one.

3.1 Belief Network for the Bayesian model

We start our model building from the LMM from Section 1.2.1, Eq. 1.3. Since we are interested in the model parameters \mathbf{E} , \mathbf{W} and σ^2 we can write the posterior using the Bayes rule:

$$f(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}) = \frac{f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\mathbf{E}, \mathbf{W}, \sigma^2)}{f(\mathbf{X})} \quad (3.1)$$

The immediate problem that we have here is how to define $f(\mathbf{E}, \mathbf{W}, \sigma^2)$. The original papers assume implicitly that $f(\mathbf{E}, \mathbf{W}, \sigma^2) = f(\mathbf{E})f(\mathbf{W})f(\sigma^2)$ which in turn allows us to write the joint probability as:

$$f(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2) = f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\mathbf{E}) f(\mathbf{W}) f(\sigma^2) \quad (3.2)$$

which in turn is equivalent with a directed acyclic graph (DAG) such as the one in Fig. 3.1a in which the Markov property is assumed [4]. Such graphs in the Bayesian framework are called Belief Networks. They are a convenient tool for describing direct influence and conditional independence assumptions between different variables. The BayesNMF-Vol model builds on these assumptions with extra hyper-parameters as depicted in Fig. 3.1b. According to graphical model theory [4], this DAG corresponds with the following joint distribution:

$$\begin{aligned} f(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2, \alpha, \beta, \gamma) \\ = f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\mathbf{E} | \gamma) f(\mathbf{W}) f(\sigma^2 | \alpha, \beta) f(\alpha) f(\beta) f(\gamma) \end{aligned} \quad (3.3)$$

Specifically, in case of BayesNMF-Vol, the parameters α , β and γ are assumed given and thus the joint distribution can further be simplified to:

$$\begin{aligned} f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2 | \alpha, \beta, \gamma) \\ = \frac{f(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2, \alpha, \beta, \gamma)}{f(\alpha, \beta, \gamma)} \\ = f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\mathbf{E} | \gamma) f(\mathbf{W}) f(\sigma^2 | \alpha, \beta) \end{aligned} \quad (3.4)$$

The last equality is due to $f(\alpha, \beta, \gamma) = f(\alpha)f(\beta)f(\gamma)$.

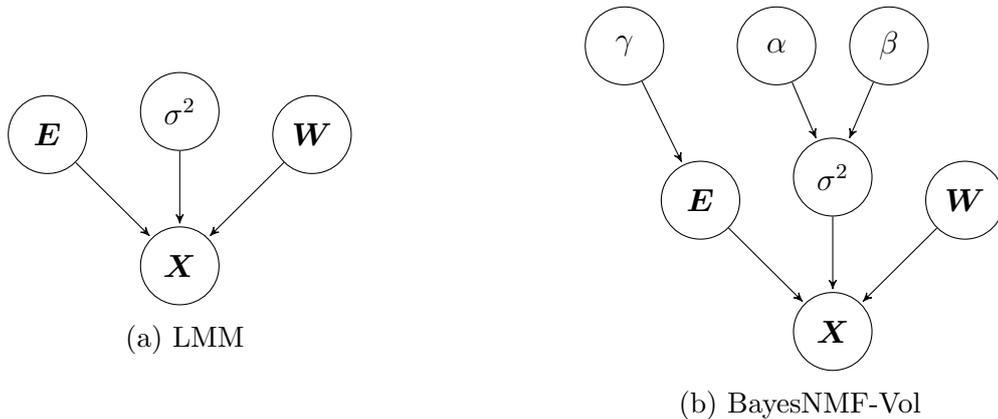


Figure 3.1: (a) Belief Network for the linear mixing model (LMM) as defined in Eq. 1.3. (b) Belief Network for the BayesNMF-Vol algorithm. Note that \mathbf{E} , \mathbf{W} and σ^2 are assumed directly independent of each other.

3.2 The partial probability density functions

We are now set to define each of the partial density functions in Eq. 3.4. The definitions are taken from Arngren [1].

3.2.1 The likelihood $f_{\mathbf{X}|\mathbf{E},\mathbf{W},\sigma^2}$

Following the LMM data model (cf. Eq. 1.3) we already derived the likelihood in Eq. 2.3. So without falling into repetition we write:

$$f(\mathbf{X}|\mathbf{E}, \mathbf{W}, \sigma^2) = \prod_{n=1}^N \prod_{b=1}^B \frac{\exp\left[-\frac{1}{2\sigma^2}(x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2\right]}{\sqrt{2\pi\sigma^2}} \quad (3.5)$$

3.2.2 The endmember prior $f_{\mathbf{E}|\gamma}$

Our first concern with the endmember prior is that it should disallow negative endmember components. Thus

$$f(\mathbf{E}) \propto \prod_{b=1}^B \prod_{m=1}^M \mathbb{I}[e_{bm} \geq 0] \quad (3.6)$$

Our second concern is that the prior should encourage the simplex spanned by the estimated endmembers to be small. Arngren's reasoning for this inclusion is because of the functional equivalence between the MAP and ML estimators as described later in section 3.3. From this, it appears that the MAP estimator would suffer the same fate as objective function L from Eq. 2.7. And thus the volume regularization is added:

$$f(\mathbf{E}|\gamma) \propto e^{-\gamma V(\mathbf{E})} \prod_b \prod_m \mathbb{I}[e_{bm} \geq 0] \quad (3.7)$$

Note that the reason for putting the volume V in the exponent is partly due to sampling reasons: the Gibbs \mathbf{E} -sampler is Gaussian in that case (cf. section 3.4.5). Furthermore, this facilitates the correspondence between the MAP and ML estimators as shown in section 3.3.

3.2.3 The abundances prior $f_{\mathbf{W}}$

The minimal requirements – as required by the LMM in Section 1.2.1 – for the abundances are the non-negativity and sum-to-one constraints. Thus:

$$f(\mathbf{W}) \propto \prod_{n=1}^N \left[\mathbb{I}[\|\mathbf{w}_n\|_1 = 1] \prod_{m=1}^M \mathbb{I}[w_{nm} \geq 0] \right] \quad (3.8)$$

Although this prior seems as least informative while incorporating the constraints, we will see later in Section 3.5.1 that it also prefers smaller endmember simplexes.

Other more complex priors could also be formulated. For example, if we acknowledge that relatively few endmembers are likely to be mixed in each pixel, then that would lead us to sparse or highly kurtotic priors for \mathbf{W} .

3.2.4 The variance prior $f_{\sigma^2|\alpha,\beta}$

The Inverse-Gamma distribution for variance σ^2 is an obvious choice since it is conjugate to the normal distribution.

$$\begin{aligned}
f(\sigma^2|\alpha, \beta) &= \mathcal{IG}(\sigma^2|\alpha, \beta) \\
&= \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp\left(-\frac{\beta}{\sigma^2}\right)
\end{aligned} \tag{3.9}$$

Arngren [1] notes that as $\alpha \rightarrow 0$ and $\beta \rightarrow 0$ the Inverse-Gamma approaches the Jeffreys prior. Now, why would one use that instead of the uniform distribution? The uniform distribution, sometimes called the Bayes-Laplace prior, has its roots in Thomas Bayes stating: ‘‘When the probability of a simple event is unknown, we may suppose all values as equally likely.’’

Jeffreys prior on the other hand has a more theoretical foundation [18]. Suppose person A and person A' are observing events from a normal distribution but use a different scale measure sp that the following equality holds: $\sigma' = c\sigma$ where $c \in \mathbb{R}_0^+$ a constant. Since their location measure is the same, they end up observing events x and x' related in the following way: $x' = c(x - \mu) + \mu$. Assuming they have the same prior information their prior probabilities are the same: $f'(\sigma, \mu) = f(\sigma, \mu)$. Assuming further independent parameters we have $f'(\sigma) = f(\sigma)$. From the change of variables theorem we can further deduce (cf. Appendix B.1) that $f'(\sigma') = f(\sigma)/c$. Combining these results gives the functional equation $cf'(c\sigma) = f(\sigma)$ which has $f(\sigma) = C^{te}/\sigma$ as the general solution. Note that this density is improper and is called the Jeffreys prior for *scale* parameters.

From a practical point of view the results do not differ much between these two priors. Nonetheless, giving more weight to smaller σ^2 values than larger ones seems more logical. And lastly, since \mathcal{IG} is conjugate to \mathcal{N} the σ^2 -sampler is \mathcal{IG} instead of a left truncated \mathcal{N} .

Now that we have defined all the constituting probability density functions of the posterior distribution, we can define the posterior.

3.2.5 The posterior $f_{\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}, \alpha, \beta, \gamma}$

We derive the posterior:

$$\begin{aligned}
&f(\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}, \alpha, \beta, \gamma) \\
&= \frac{f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2 | \alpha, \beta, \gamma)}{f(\mathbf{X} | \alpha, \beta, \gamma)} \\
&= \frac{f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\mathbf{E} | \gamma) f(\mathbf{W}) f(\sigma^2 | \alpha, \beta)}{\int_{\mathbf{E}} \int_{\mathbf{W}} \int_{\sigma^2} f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\mathbf{E} | \gamma) f(\mathbf{W}) f(\sigma^2 | \alpha, \beta) d\sigma^2 d\mathbf{W} d\mathbf{E}} \\
&\propto f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\mathbf{E} | \gamma) f(\mathbf{W}) f(\sigma^2 | \alpha, \beta) \\
&\propto \prod_{n=1}^N \prod_{b=1}^B \frac{\exp\left[-\frac{1}{2\sigma^2} (x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2\right]}{\sqrt{2\pi\sigma^2}} \\
&\quad \times e^{-\gamma V(\mathbf{E})} \prod_{b=1}^B \prod_{m=1}^M \mathbb{I}[e_{bm} \geq 0] \\
&\quad \times \prod_{n=1}^N \left[\mathbb{I}[\|\mathbf{w}_n\|_1 = 1] \prod_{m=1}^M \mathbb{I}[w_{nm} \geq 0] \right] \\
&\quad \times \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp\left(-\frac{\beta}{\sigma^2}\right)
\end{aligned} \tag{3.10}$$

Note that the integral in the denominator is just a constant.

3.3 Relation between MAP and ICE estimators

Several authors [3], [8] have noted an interesting connection between the geometrical approaches such as the ICE algorithm and the Bayesian MAP estimator:

$$\begin{aligned}
(\hat{\mathbf{E}}, \hat{\mathbf{W}})_{MAP} &= \arg \max_{\mathbf{E}, \mathbf{W}} f_{\mathbf{E}, \mathbf{W}}(\mathbf{E}, \mathbf{W} | \mathbf{X}, \sigma^2, \alpha, \beta, \gamma) \\
&= \arg \max_{\mathbf{E}, \mathbf{W}} f_{\mathbf{X}}(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f_{\mathbf{E}}(\mathbf{E} | \gamma) f_{\mathbf{W}}(\mathbf{W}) \\
&= \arg \min_{\mathbf{E}, \mathbf{W}} -\log f_{\mathbf{X}}(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) - \log f_{\mathbf{E}}(\mathbf{E} | \gamma) - \log f_{\mathbf{W}}(\mathbf{W}) \\
&= \arg \min_{\mathbf{E}, \mathbf{W}} \|\mathbf{X} - \mathbf{W} \mathbf{E}^T\|_F^2 + \gamma V + C^{te}
\end{aligned} \tag{3.11}$$

which is similar in form to L_{reg} in Eq. 2.7. In the ICE algorithm, the estimates are obtained by minimizing a two-term objective function L_{reg} where $-\log f_{\mathbf{X}}(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2)$ plays the role of a data fitting criterion and $-\log f_{\mathbf{E}}(\mathbf{E} | \gamma) - \log f_{\mathbf{W}}(\mathbf{W})$ of a penalization. Conversely, from a Bayesian perspective, assigning prior distributions to \mathbf{E} and \mathbf{W} is a convenient way to ensure physical constraints inherent to the observation model. Thus such geometrical algorithms have statistical foundation as Bayesian MAP estimators.

3.4 The BayesNMF-Vol sampler algorithm

3.4.1 Choice of a sampler

Sampling from the posterior is not straightforward. The first difficulty is that we do not know the exact distribution since the denominator is not known. The Metropolis-Hastings (MH) algorithm does not require the denominator, and thus is a natural choice. Problem is the amount of parameters. The MH algorithm is highly dependent on a good proposal function, without which the acceptance rate is halved with each extra dimension. A viable solution is a Gibbs sampling approach which allows sampling at univariate level – at the expense of meticulous calculations.

3.4.2 The Gibbs sampler

The general Gibbs sampling procedure is described in [9]. We use it here to sample $(\mathbf{E}, \mathbf{W}, \sigma^2)$ from the posterior $f_{\mathbf{E}, \mathbf{W}, \sigma^2 | \mathbf{X}, \alpha, \beta, \gamma}$. The specific procedure is defined in Algorithm 2. The algorithm generates a Gibbs sequence of random variables $((\mathbf{E}^{(0)}, \mathbf{W}^{(0)}, \sigma^{2(0)}), \dots, (\mathbf{E}^{(k)}, \mathbf{W}^{(k)}, \sigma^{2(k)}))$ that constitute a Markov chain for which the stationary distribution is the posterior in which we are interested.

Algorithm 2 (BayesNMF-Vol sampling algorithm)

Initialize $\mathbf{E}^{(0)}$ and $\mathbf{W}^{(0)}$ in a random feasible way.

for $i = 1, 2, \dots$ to k **do**

$\sigma_1^{2(i)} \sim f_{\sigma^2}(\sigma^2 | \mathbf{E}^{(i-1)}, \mathbf{W}^{(i-1)}, \mathbf{X}, \alpha, \beta, \gamma)$

$w_{11}^{(i)} \sim f_{\mathbf{W}}(w_{11} | \sigma_1^{2(i)}, \mathbf{E}^{(i-1)}, w_{12}^{(i-1)}, \dots, w_{NM}^{(i-1)}, \mathbf{X}, \alpha, \beta, \gamma)$

\vdots

$w_{NM}^{(i)} \sim f_{\mathbf{W}}(w_{NM} | \sigma_1^{2(i)}, \mathbf{E}^{(i-1)}, w_{11}^{(i)}, \dots, w_{NM-1}^{(i)}, \mathbf{X}, \alpha, \beta, \gamma)$

$e_{11}^{(i)} \sim f_{\mathbf{E}}(e_{11} | \sigma_1^{2(i)}, \mathbf{W}^{(i)}, e_{12}^{(i-1)}, \dots, e_{BM}^{(i-1)}, \mathbf{X}, \alpha, \beta, \gamma)$

\vdots

$e_{BM}^{(i)} \sim f_{\mathbf{E}}(e_{BM} | \sigma_1^{2(i)}, \mathbf{W}^{(i)}, e_{11}^{(i)}, \dots, e_{BM-1}^{(i)}, \mathbf{X}, \alpha, \beta, \gamma)$

end for

return $(\mathbf{E}, \mathbf{W}, \sigma^2)$ as Gibbs sample sequence.

Advantages to the MH algorithm are that we do not need to tune the proposal distribution and we have no inefficiency due to rejected proposals. The downside is that the Gibbs-progress can be stalled by highly correlated parameters.

In the following sections we derive each of the conditional distributions needed by the Gibbs sampler and define appropriate ways of sampling. These derivations are based on the technical paper [2]. In case our own derivations are different or include omitted steps, we put them into the Appendix.

3.4.3 Sampling σ^2

Because we have chosen a conjugate prior for the noise variance, its conditional distribution has the same functional form as the prior: An Inverse-Gamma.

$$\begin{aligned}
 f(\sigma^2 | \mathbf{E}, \mathbf{W}, \mathbf{X}, \alpha, \beta, \gamma) &= \frac{f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2 | \alpha, \beta, \gamma)}{\int_{\sigma^2} f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2 | \alpha, \beta, \gamma) d\sigma^2} \\
 &= \frac{f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\sigma^2 | \alpha, \beta)}{\int_{\sigma^2} f(\mathbf{X} | \mathbf{E}, \mathbf{W}, \sigma^2) f(\sigma^2 | \alpha, \beta) d\sigma^2} \\
 &= f(\sigma^2 | \mathbf{E}, \mathbf{W}, \mathbf{X}, \alpha, \beta) \\
 &= \mathcal{IG}(\sigma^2 | \bar{\alpha}, \bar{\beta})
 \end{aligned} \tag{3.12}$$

Note that in the second equation $f(\mathbf{E} | \gamma)$ and $f(\mathbf{W})$ can be removed from the integral since they are just constants. Thus $f_{\sigma^2 | \mathbf{E}, \mathbf{W}, \mathbf{X}, \alpha, \beta}$ is independent of γ given the other parameters. The Inverse-Gamma parameters, as calculated in [2], are given by

$$\bar{\alpha} = \alpha + \frac{1}{2}NB \tag{3.13}$$

$$\bar{\beta} = \beta + \frac{1}{2}\|\mathbf{X} - \mathbf{W}\mathbf{E}^T\|_F^2 \tag{3.14}$$

Parameter $\bar{\alpha}$ controls the shape and is *constant*. The larger it is, the more the density will shift to 0 and the less variable the outcome. $\bar{\beta}$ on the other hand controls the scale and is directly proportional to the error. The more pixels lie out of the simplex, the larger it is, and thus the larger the variance.

Note that $\bar{\alpha}$ and $\bar{\beta}$ most likely will never be equal to zero even if we set $\alpha = 0$ and $\beta = 0$ to simulate Jeffreys prior.

3.4.4 Sampling \mathbf{W}

$$\begin{aligned}
 f(\mathbf{W}|\mathbf{E}, \sigma^2, \mathbf{X}) &= \prod_{n=1}^N f(\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n) \\
 &\propto \prod_{n=1}^N \left[\mathcal{N}(\mathbf{w}_n|\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_w) \mathbb{I}[\|\mathbf{w}_n\|_1 = 1] \prod_{m=1}^M \mathbb{I}[w_{nm} \geq 0] \right]
 \end{aligned} \tag{3.15}$$

where the parameters are given by

$$\boldsymbol{\Sigma}_w^{-1} = \frac{\mathbf{E}^T \mathbf{E}}{\sigma^2} \tag{3.16}$$

$$\boldsymbol{\mu}_n = (\mathbf{E}^T \mathbf{E})^\dagger \mathbf{E}^T \mathbf{x}_n \tag{3.17}$$

The derivation is worked out in Appendix B.3. First note that we can concentrate on $f_{\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n}$ instead of the whole $f_{\mathbf{W}|\mathbf{E}, \sigma^2, \mathbf{X}}$. Next, note that $f_{\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n}$ is multivariate Gaussian, constrained on a unit-simplex. Fig. 3.2 shows how this probability density might look like for $\mathbf{w}_n \in \mathbb{R}^2$. Note that in this case, the feasible region is a straight line since the unit-simplex is 1-dimensional.

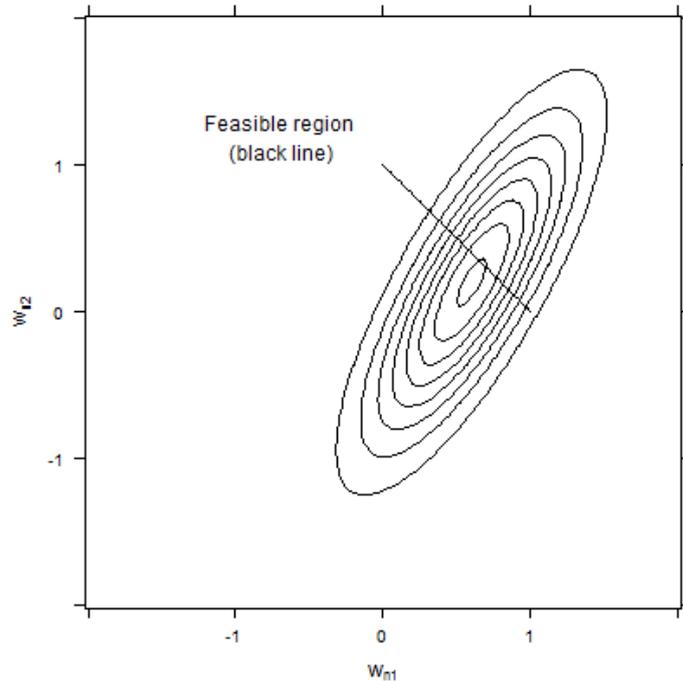


Figure 3.2: The contours are from $\mathcal{N}(\mathbf{w}_n|\boldsymbol{\mu}_n, \boldsymbol{\Sigma})$ with $\boldsymbol{\Sigma} = [0.2, 0.25; 0.25, 0.5]$, $\boldsymbol{\mu}_n = (0.8, 0.2)$ and $\mathbf{w}_n \in \mathbb{R}^2$. The support of the density $f_{\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n}$ is only along the black line.

Sampling from $f_{\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n}$ at this stage can be easily done with a MH sampler, using a uniform Dirichlet proposal. Dirichlet random variables meet the non-negativity as well

as the sum-to-one constraints and thus the only thing which remains to be checked is compliance to the Normal density. This way of sampling is very efficient and feasible for relatively low M .

Removing equality constraint from $f_{w_n|\mathbf{E},\sigma^2,\mathbf{x}_n}$

A more complex approach is a univariate Gibbs sampler which works for any M . Note that we can not condition on $f_{w_n|\mathbf{E},\sigma^2,\mathbf{x}_n}$ and go directly for w_{nm} since the sum-to-one constraint is deterministic, implying $\text{Var}(w_{nm}|w_{n\tilde{m}}) = 0$. Our approach is inspired by Schmidt [23] who describes a general approach of how sampling can be done from a distribution constrained with any number of equality and inequality constraints. The idea is to first map the distribution to an affine subspace where the equality constraints hold. In that subspace the distribution is only truncated (because the inequality constraints still hold), but not degenerate.

The original papers refer to this novel idea, but do not translate it to the problem at hand. In this Section we derive the required components while relying on this general idea.¹ Thus we sample \mathbf{y} from a reduced but equivalent space through the relationship $w_n = \mathbf{A}\mathbf{y} + \mathbf{b}$ using

$$f(\mathbf{y}|\mathbf{E}, \sigma^2, \mathbf{x}_n) = \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \mathbb{I}[\|\mathbf{y}\|_1 \leq 1] \prod_{m=1}^{M-1} \mathbb{I}[\mathbf{y} \geq 0] \quad (3.18)$$

where

$$\boldsymbol{\Sigma}_y^{-1} = \frac{\mathbf{A}^T \mathbf{E}^T \mathbf{E} \mathbf{A}}{\sigma^2} \quad (3.19)$$

$$\boldsymbol{\mu}_y = (\mathbf{A}^T \mathbf{E}^T \mathbf{E} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{E}^T [\mathbf{x}_n - \mathbf{E}\mathbf{b}] \quad (3.20)$$

The specification of \mathbf{A} and \mathbf{b} with full derivation of these results is given in Appendix B.4. In this new \mathbf{y} -space, the equality constraint does not hold and the density is a mere truncated Gaussian distribution in which the covariance matrix is much more stable. We can sample from this distribution using a Gibbs sampler for which the univariate density $f_{y_k|\mathbf{y}_{\bar{k}},\mathbf{E},\sigma^2,\mathbf{x}_n}$ is derived in Appendix B.5. The downside of this sampler compared to MH is slower sampling. The upside is that nothing ever gets rejected which results in a visible improvement in impure endmember detection.

3.4.5 Sampling \mathbf{E}

The conditional distribution of the endmembers \mathbf{E} arises from the product of the Gaussian likelihood and the volume penalizing prior.

¹We do not base ourselves on the results of Schmidt [23] – only his general idea – because we are not sure that his resulting inequality constraints match with those that follow from our own derivations. The general procedure does not seem to account completely for additional inequality constraints that arise once the mapping is done. This requires further research.

$$\begin{aligned}
f(\mathbf{E}|\mathbf{W}, \sigma^2, \mathbf{X}, \gamma) &= \prod_{b=1}^B f(\mathbf{e}_b|\mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma) \\
&\propto \prod_{b=1}^B \left[\mathcal{N}(\mathbf{e}_b|\boldsymbol{\mu}_b, \boldsymbol{\Sigma}_E) \prod_{m=1}^M \mathbb{I}[e_{bm} \geq 0] \right]
\end{aligned} \tag{3.21}$$

where the parameters are given by

$$\boldsymbol{\Sigma}_E^{-1} = \frac{\mathbf{W}^T \mathbf{W}}{\sigma^2} + 2\gamma \left(\frac{\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M}{M-1} \right) \tag{3.22}$$

$$\boldsymbol{\mu}_b = \frac{\boldsymbol{\Sigma}_E \mathbf{W}^T \mathbf{x}_b}{\sigma^2} \tag{3.23}$$

The derivation of the parameters is given in Appendix B.6. We see that $f_{\mathbf{e}_b|\mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma}$ is a multivariate Normal distribution constrained to the positive domain. The density $f_{e_{bm}|e_{b\bar{m}}, \mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma}$ for use within a Gibbs sampler is derived in Appendix B.7. The results are different and computationally more efficient, but equivalent to the results in [2].

Concerning the inverse, we have already touched upon a similar issue in Section 2.3.4. We conclude that it is highly unlikely that $\boldsymbol{\Sigma}_E$ is singular.

3.5 Model peculiarities

3.5.1 Intrinsic regularization of the Bayesian approach

An interesting feat is that the Bayesian approach is inherently regularized. Fig. 3.3 shows that even if $\gamma = 0$, the endmember simplex will be stable and never much larger than the data points. Arngren already noticed this in [1] but had no explanation. Two years later, in [3], he turns up with a mathematical explanation for one dimension. The main insight is that we are not so much interested in $(\hat{\mathbf{E}}, \hat{\mathbf{W}})_{MAP}$ from Eq. 3.11 but in $(\hat{e}_{bm})_{MAP}$ and $(\hat{w}_{nm})_{MAP} \forall n, m, b$. Thus we are interested in marginal distributions. While the MAP estimator of the joint distribution is ill-defined, the MAP estimators of the components are not! And since we are dealing with indicator functions in priors the Riemann integration cannot be used, and one must resort to Lebesgue integration to come up with a closed-form solution. Arngren [3] does this for one dimension.

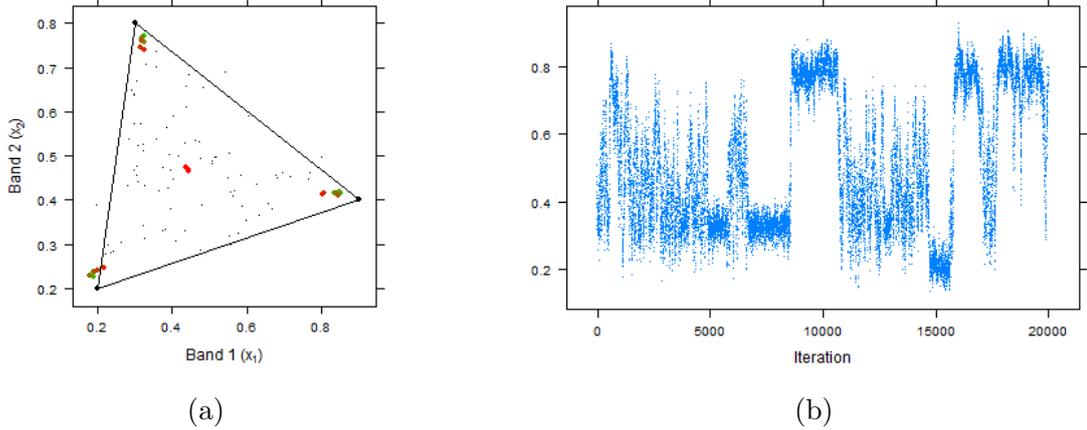


Figure 3.3: Data for (a) is simulated as a 2-band 9×9 image with $\text{SNR}_{\text{dB}} = 15$ dB. The MAP endmember estimates are shown from green to red for $\gamma = \{0, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 140, 1000\}$. The estimates are based on 3000 samples with 500 burn-in. The estimates are stable until we reach $\gamma = 140$. For $\gamma = 150$ and above the endmember approximations become multi-modal, as the trace in (b) shows for e_{11} . Once the simplex reaches minimal size, the MAP estimators become stable again.

We can understand this intrinsic regularization conceptually by looking at the prior for \mathbf{W} in Eq. 3.8. The key for understanding is seeing that this prior conveys more information than just the non-negativity and sum-to-one constraints. It also says that each w_{nm} that satisfies the constraints is equally likely. In Fig. 3.4 we have two endmember simplexes for the same dataset. Given \mathbf{E} , for the left figure the abundances w_{nm} are uniformly distributed. But for the right figure, all the abundances w_{nm} are around 0.35. Our prior on \mathbf{W} will thus favor the left simplex.

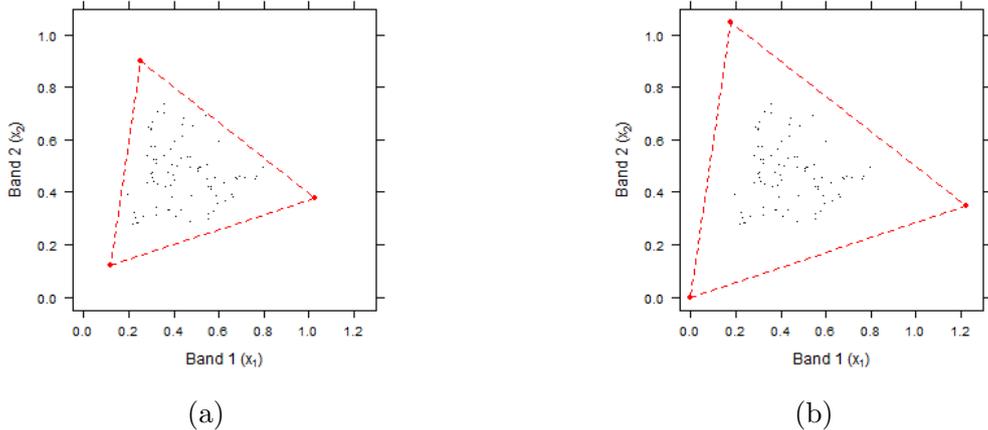


Figure 3.4: Simplex in (a) is the MAP estimate for \mathbf{E} with $\gamma = 0$ such that all w_{nm} are equally likely. Simplex in (b) on the other hand is such that most w_{nm} fall within the 0.3-0.4 range.

3.6 Conclusion

In this chapter we defined the Bayesian framework for the LMM model using a DAG with the Markov property. Under these assumptions we defined the priors for \mathbf{E} , \mathbf{W} and σ^2 as uninformative as possible, while fulfilling the requirements of the LMM. For \mathbf{E} and \mathbf{W} this resulted in improper priors only asserting non-negativity and sum-to-zero constraints through indicator functions. A volume regularization on \mathbf{E} was also made possible through the \mathbf{E} -prior, controlled by the hyper-parameter γ . For σ^2 we used the Jeffreys prior for scale parameters. We derived it and showed how it can be implemented as the limiting case of the Inverse-Gamma distribution. Furthermore, we showed how the MAP estimators of this model are related to the ICE estimators. We further went into greater detail than the original papers in deriving all the necessary conditional distributions for the Gibbs sampling procedure. Notable here is the sampling procedure from a degenerate distribution described in Section 3.4.4. We first mapped this distribution to a subspace in which the equality constraints hold, and sampled from that subspace. Then we back-transformed the samples to the original space.

The simulation results from Section 3.5.1 indicate that the Bayesian model is intrinsically volume-regularized and thus that γ can safely be set to zero. We explained the reason for this conceptually and referred to some mathematical proofs in other papers. We also noted that although the MAP estimator is ill-defined for the whole distribution, it is not once we consider the marginals $(\hat{e}_{bm})_{MAP}$ and $(\hat{w}_{nm})_{MAP}$. It is not always clear from these papers if, and how, $(\hat{\mathbf{E}}, \hat{\mathbf{W}})_{MAP}$ are computed. In subsequent chapters we will use the marginal MAP estimators. Let us now turn to some real data examples.

Chapter 4

Synthetic data

4.1 Image generation procedure

For simulating hyperspectral images we have three parameters to play with: the number of pixels N , the number of endmembers M and the number of bands B . Also, the image content (i.e. the ordering of the pixels) can be adjusted.

We confine ourselves to the 50×50 -pixel image ($N = 2500$) in Fig. 4.1. For representational clarity we use only $B = 2$ bands and $M = 3$ synthetic endmembers as shown in Fig. 4.3d. The 3 abundance maps shown in Fig. 4.3a-c were generated using Markov Random Fields (MRF) with Matern covariance function. Once mixed with the endmembers, they result in Fig. 4.1a. Then we add $\text{SNR}_{\text{dB}} = 15$ dB Gaussian noise which results in Fig. 4.1b.¹

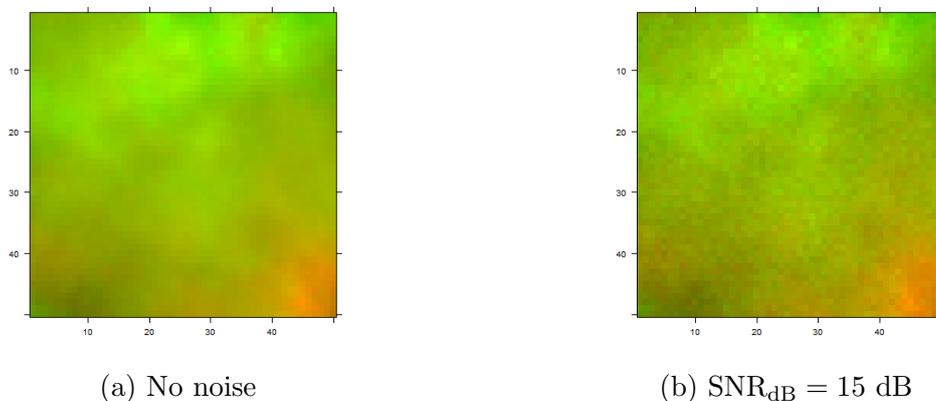


Figure 4.1: False color synthetic image with (a) and without (b) noise. The two bands are represented as red and green.

In the following section we run the ICE, ICE-S and BayesNMF-Vol algorithms on the synthetic image and compare the estimated endmembers and abundance maps with the originals. Then we select an appropriate error metric for the endmembers. And in the subsequent section we explore the whole parameter space of the three models to find the optimal parameters.

¹Clark [10] reports AVIRIS data exceeding SNR of 500 at most wavelengths. This corresponds with $\text{SNR}_{\text{dB}} = 10 \log_{10} 500 \sim 27$ dB for real data.

4.2 Analysis

ICE, ICE-S and BayesNMF-Vol algorithms were run to extract the endmembers and abundance maps. For ICE and ICE-S we use the threshold 0.9999 of the ratio between the current and the previous outcome of objective function as the stopping criterion. For BayesNMF-Vol algorithm we use 3000 samples with 500 samples as burn-in. To make the comparison between BayesNMF-Vol and ICE algorithms meaningful to some degree, we use the MAP point estimators in the Bayesian setting. We extract them from the MC samples using the kernel density estimation technique with Gaussian kernel as implemented in the R *density()* function.

The resulting endmembers and abundance maps are shown in Fig. 4.3 for comparison. The image can be split in three groups. The first group is the first row which consists of original abundance maps and endmembers. The three subsequent rows comprise the second group consisting of ICE, ICE-S and BayesNMF-Vol estimates with suggested parameters. For the ICE algorithm we use $\mu = 0.01$ as suggested by Berman [7]. For ICE-S we chose 98 % spatial regularization with same volume regularization as ICE, mainly to contrast the smoothing effect with ICE. For BayesNMF-Vol we set the parameter $\gamma = 0$ to show the intrinsic regularization. The last group consists of two last rows of ICE and ICE-S estimates with optimal parameters which are found in the next section.

Overall we see similar endmember estimates for ICE and ICE-S algorithms. As long as the volume regularization is kept the same, spatial regularization can be used to smoothen the abundance maps without changing the endmember estimates much.

4.3 Evaluation

Now that we have the endmembers and abundance maps, we want an objective means to compare them with the real ones as a way of evaluating the algorithms. For that we first need to select an error measure.

4.3.1 Error metric selection

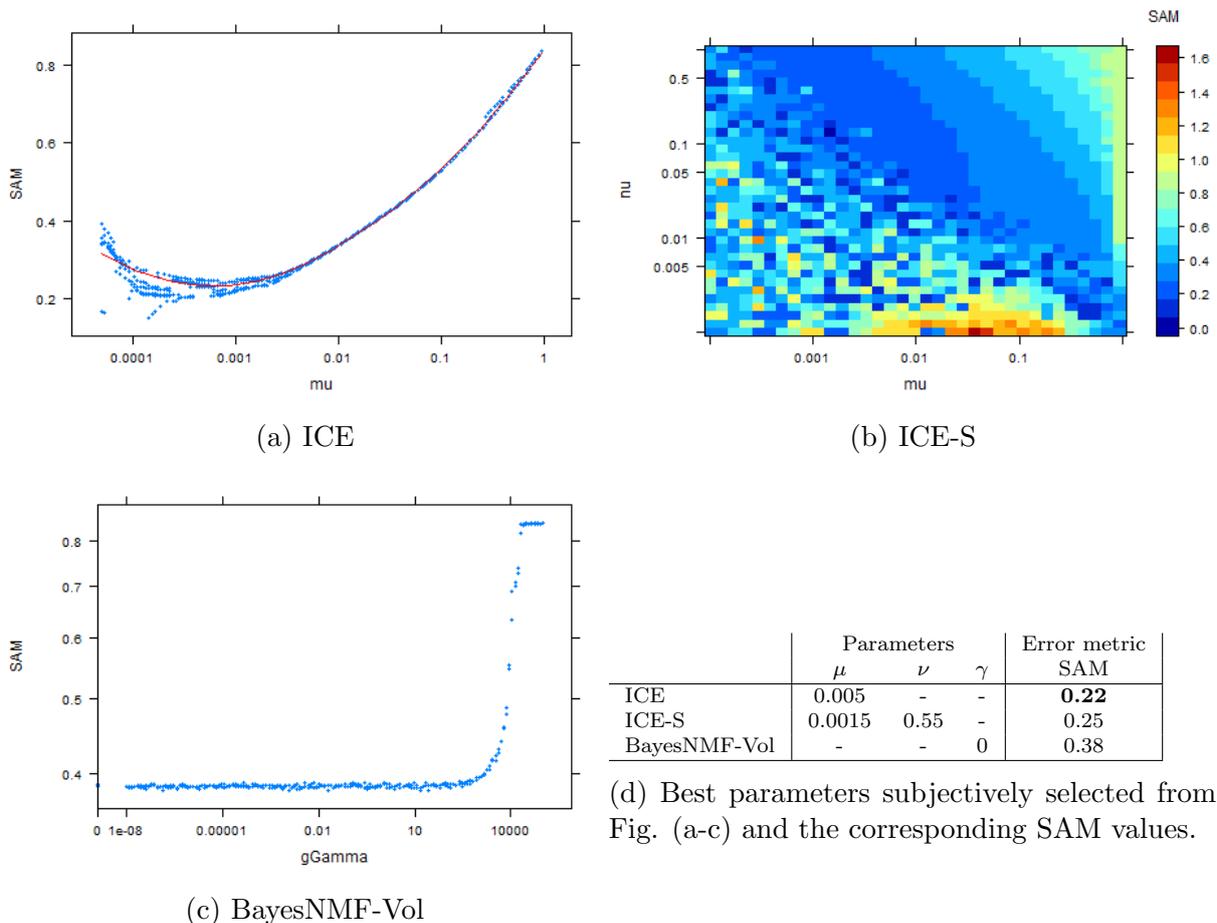
We are mainly interested in endmember identification here. So an obvious measure candidate is the mean squared error (MSE). The problem with this measure is that it is dependent on variations in albedo (i.e. the brightness or gain factor) of the spectra [15]. Variations in albedo can occur due to different topographic illumination effects, different device calibration, etc. A more appropriate measure would be one that is only dependent on the shape of the spectra, but not on albedo. Such a measure is the spectral angle mapper (SAM) [19]. This method quantifies the similarity between spectra by calculating the "angle" between them using the following formula:

$$\alpha_{SAM} = \cos^{-1} \left[\frac{\mathbf{e}_{ref} \mathbf{e}_{test}}{\|\mathbf{e}_{ref}\| \|\mathbf{e}_{test}\|} \right] \quad (4.1)$$

This measure will not change for \mathbf{e}_{ref} and \mathbf{e}_{test} when \mathbf{e}_{test} is multiplied with some $c \in \mathbb{R}_0^+$ and thus the measure is insensitive to albedo. Further more, when $\mathbf{e}_{ref} = \mathbf{e}_{test}$ – and for any $c\mathbf{e}_{test}$ – this measure will result in 0. For a more in-depth discussion of this measure, and a comparison with a variant of a MSE-based measure, one can consult Dennison [15].

4.3.2 Results

Fig. 4.2 gives an overview of SAM values in function of the parameters μ , ν and γ for ICE, ICE-S and BayesNMF-Vol algorithms respectively on the above mentioned synthetic image. A few noteworthy things are visible in these graphs. The ICE algorithm seems to have only one global minimum for μ down to $\mu = 0.005$. Thereafter seem to be lots of local minima in which the algorithm gets stuck. The optimal and least variable μ seems to be around 0.005. For ICE-S, the more noise we see in the Fig. 4.2b heat map, the more local minima there are in that region. Best parameters are the ones with least local minima, and the smallest SAM value. We estimate that subjectively to be around $\mu = 0.0015$ and $\nu = 0.55$. For the BayesNMF-Vol, the parameter γ does not play a mayor role in interval $]-\inf, 100]$.



(d) Best parameters subjectively selected from Fig. (a-c) and the corresponding SAM values.

Figure 4.2: Figures (a-c) show SAM values in function of the model hyper-parameters. The ICE algorithm was run 5×150 times for different values of μ . The ICE-S algorithm was run 35×35 times for different values of μ and ν . The BayesNMF-vol algorithm was run 5×150 times for different values of γ . At each run the initialization was different so that in case of the ICE algorithm, the minima differ. The red curve in (a) is a quadratic LOESS curve. The runs are based on the aforementioned synthetic image.

Based on these SAM values we conclude that ICE and ICE-S perform similarly. For the used parametrization, ICE gives us the best SAM value of 0.22. Note though that SAM values can go down as far as 0.03 for ICE-S, but these minima all seem local. There are also seemingly stable regions where SAM of 0.20 is attained. But due to large parameter

space and the slowness of the algorithm it is difficult to pinpoint them. The generation of Fig. 4.2b which consists of 35×35 points took 10 days on two physical processor cores.

4.4 Conclusion

We have constructed a synthetic image based on 2 bands and 3 endmembers. We have further run ICE, ICE-S and BayesNMF-Vol on the synthetic image using a wide range of hyper-parameter configurations. We have compared the estimated endmembers with the original ones using the SAM metric. We conclude that although ICE-S was able to produce the least SAM errors, they all seemed to be local minima. In the SAM range around 0.20-0.30 both ICE and ICE-S were stable and performed similarly. Spatial regularization of ICE-S can be used to smoothen the abundance maps while keeping similar endmember estimates as of ICE, as long as the same volume regularization is used for both algorithms.

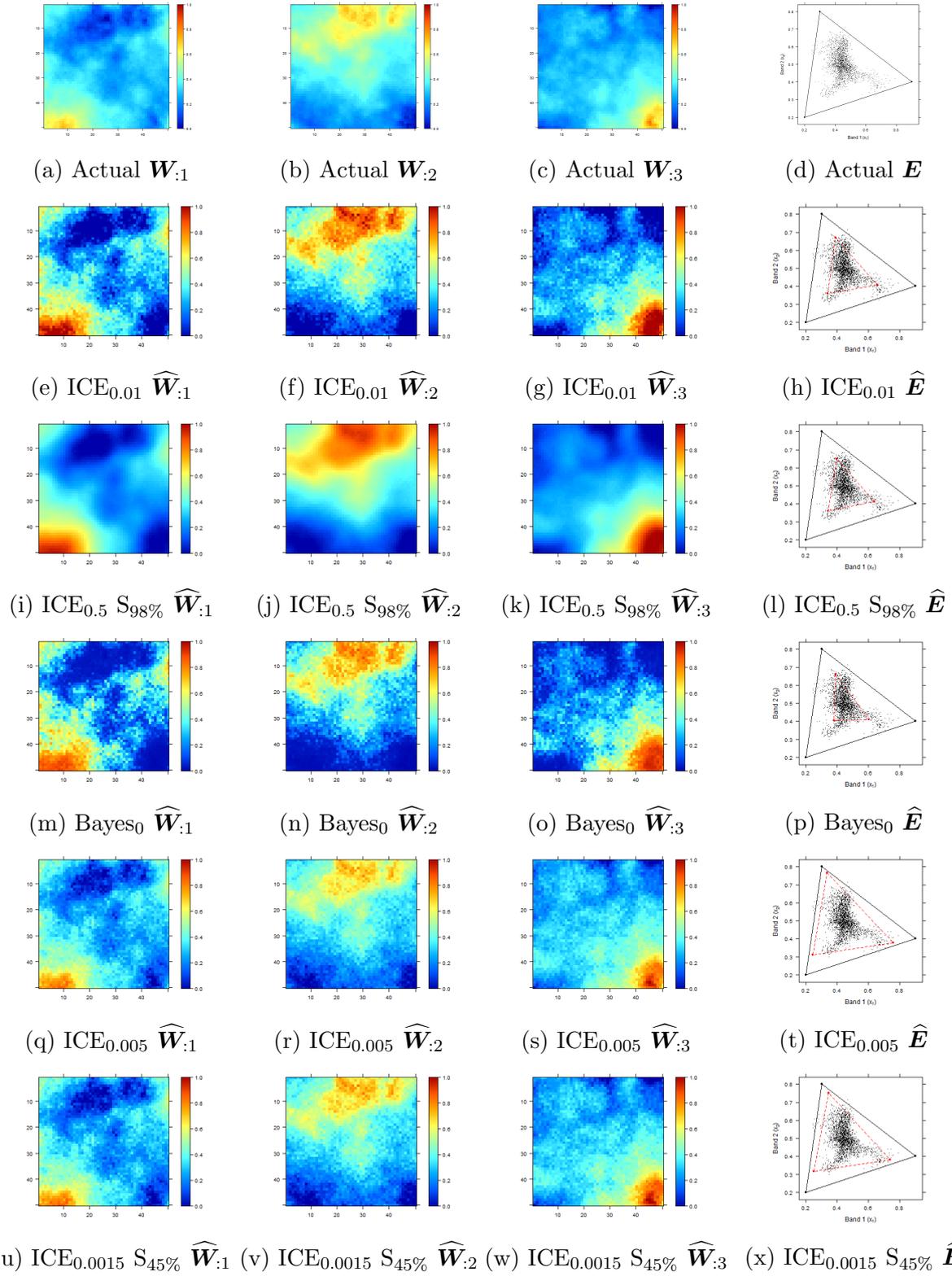


Figure 4.3: Endmembers and abundance maps of the synthetic image for ICE, ICE-S and BayesNMF-Vol. The subscript for ICE is the μ hyper-parameter and the percentage for S is related to the $1 - \nu$ hyper-parameter. The Bayes subscript is the γ hyper-parameter.

Chapter 5

Real AVIRIS data - Cuprite

This chapter illustrates the discussed algorithms on a real hyperspectral dataset. This real dataset was acquired over Cuprite Hills (Navada, USA) in 1995 by the NASA/JPL “Airborne Visual and Infra-Red Imaging Spectrometer” (AVIRIS). The instrument was flown in an ER-2 aircraft (a modified U-2 spy plane) at 20 km. The spacing between pixels is 15 m, and the size of each pixel is about 18 m. With 224 bands, a spectral sampling and bandpass of ~ 10 nm and signal-to-noise ratios exceeding 500 at most wavelengths, AVIRIS data are widely used for terrestrial remote sensing. This dataset is freely available [14].

In the first section we explain hyperspectral data preprocessing and apply it to the AVIRIS data. The second section is devoted to the selection of the appropriate “alunite hill” scene for data analysis. In subsequent Section 5.3 we run the ICE, ICE-S and BayesNMF-Vol algorithms on the selected scene. We summarize our results in Section 5.4.

5.1 Preprocessing AVIRIS data

This Section is mainly based on the *Preprocessing AVIRIS Data Tutorial*¹. The preprocessing of the data can be seen as a 3 step process.

5.1.1 Raw data

First we have the raw quantized pixel data as acquired by the sensor. NASA/JPL always processes the AVIRIS data to remove geometric and radiometric errors associated with the motion of the aircraft used during data collection.

5.1.2 Radiance data

Radiance is the amount of radiation coming from an area. Due to the variance in band-sensitivity of the sensor, the raw data has to be calibrated to compensate this variability. This is done by multiplying the raw data by a series of gain values, one for each band resulting in $W/(cm^2 * sr * nm)$.

Radiance image includes the radiation effects of the sun. In fact, the spectrum of a radiance pixel closely matches that of the solar irradiance curve, i.e. the solar spectrum.

¹<http://www.harrisgeospatial.com/portals/0/pdfs/envi/PreprocessAVIRIS.pdf>

Besides, atmospheric gas absorptions also cause specific irregularities in the spectrum, as shown in Fig. 5.1.

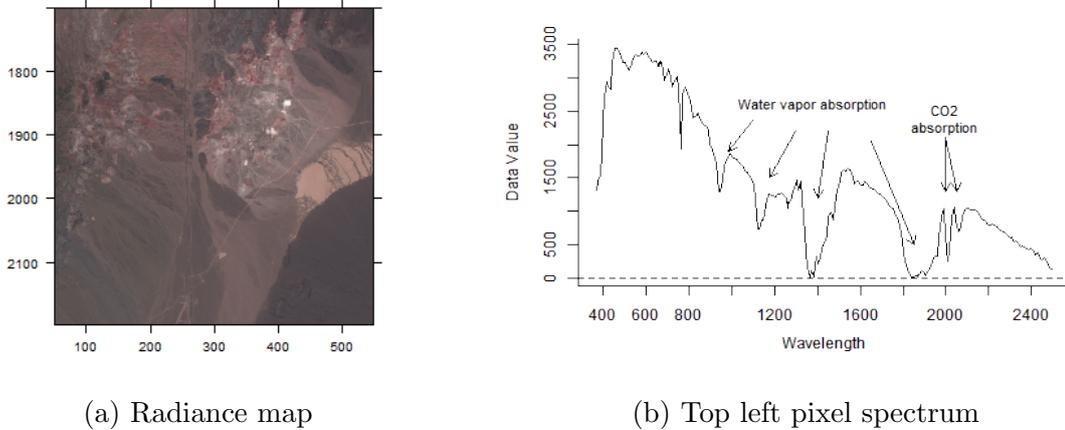


Figure 5.1: Left (a) is the false color Cuprite radiance image as obtained from [14]. The x and y axis denote the coordinates of each pixel from the original ENVI data which is much larger. For the RGB colors, the wavelengths 753.1287 nm, 530.8180 nm and 482.1898 nm were taken with intensities of each band normalized with the maximal band value. The right image (b) is the radiance spectrum of the top left pixel. The shape of the solar irradiance is clearly visible, and so are the water vapor and CO_2 absorptions.

For hyperspectral data analysis, one removes the effects of solar irradiance and atmosphere by calibrating the data to reflectance.

5.1.3 Reflectance data

Reflectance is the proportion of the radiation reflected off a surface to the radiation striking it. In hyperspectral data analysis, materials are identified by their reflectance spectra. So calibrating the data to reflectance is an important step towards identifying materials from an image. An atmospheric correction tool can remove the effects of atmospheric scattering and gas absorptions, to produce reflectance data as shown in Fig. 5.2.

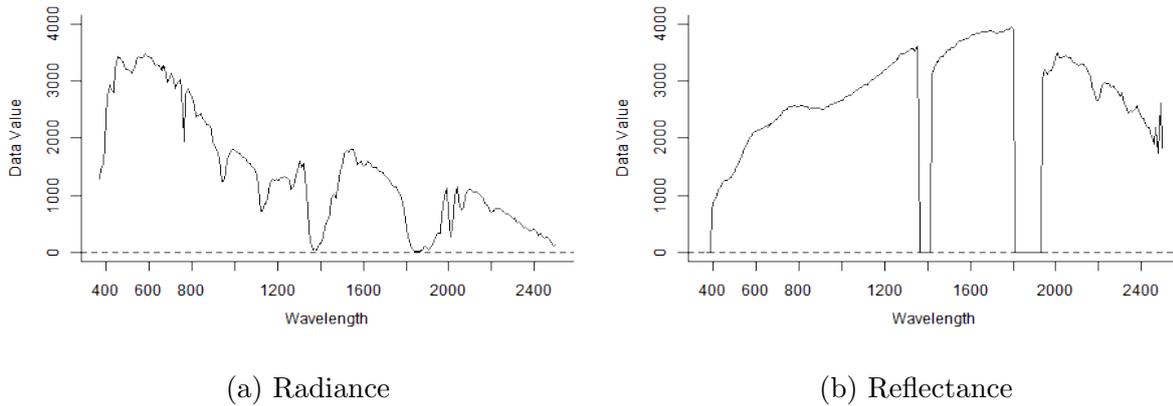


Figure 5.2: The left image (a) is the radiance spectrum of the top-left pixel (cf. Fig. 5.1b). The right image (b) shows the same spectrum processed by a atmospheric correction tool to produce the reflectance spectrum.

Note that the imagery may still have variations in illumination due to topography after this step.

5.1.4 Data cleanup

There can be some irregularity in the reflectance data caused by the atmospheric correction tool, mostly around the absorption bands. Fig. 5.3 depicts the intensity distribution in function of the band which clearly shows outlying values around the absorption wavelengths. 27 bands in total have been removed due to negative, outlying or noisy values, resulting in 197 remaining bands.

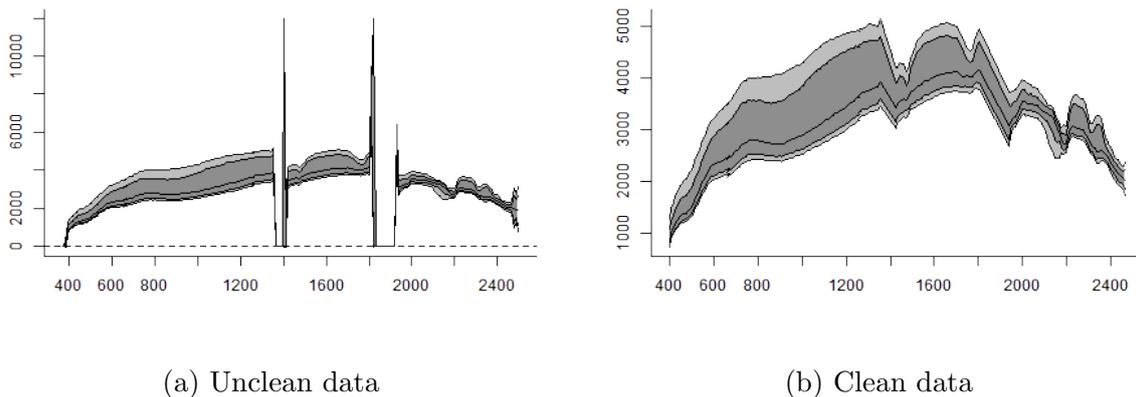


Figure 5.3: Plot (a) shows the 5th, 16th, 50th, 84th and 95th percentile per band of the uncleaned data. Plot (b) shows the same data, but with bands 366 – 385 nm, 1363 – 1413 nm and 1821 – 1918 nm removed due to negative or zero values and bands 1811 nm and 1928 nm removed due to extreme outlying values. The last tree bands 2477 – 2497 nm have also been removed due to seemingly high noise (cf. [20]).

5.2 Scene selection

The three algorithms have been evaluated on the “alunite hill” that appears in the Cuprite scene as depicted in Fig. 5.4. The geological properties of this area have been extensively investigated by geologists in [10], and [11]. The results of these studies are summarized in Fig. 5.5. The endmembers are dominated by three materials: muscovite, alunite, and kaolinite. They allow us to assess, at least subjectively, the accuracy of the three models.

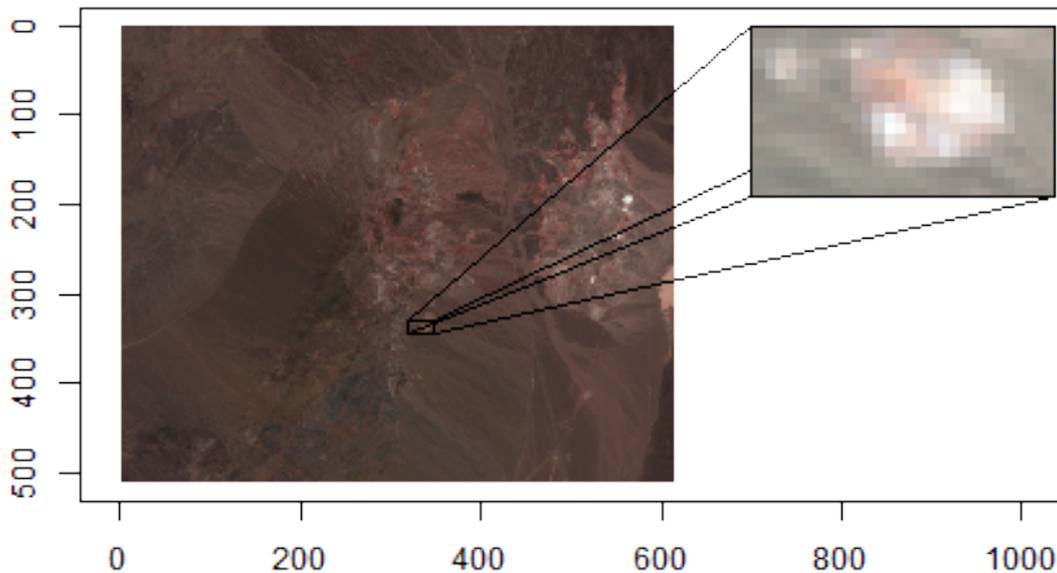


Figure 5.4: Left is the false color Cuprite reflectance image as obtained from [14]. Right is the zoomed in area of the “alunite hill” (16 x 28 pixels) used for material analysis. False colors are generated as described in Fig. 5.2.

5.3 Analysis

The 16 x 28 subimage with 197 bands, depicted in Fig. 5.4, has been unmixed using the ICE, ICE-S and BayesNMF-Vol algorithms using the same parameters as in the previous chapter. To assess the performance of the algorithms we need reference endmembers and abundance maps.

We compare the predicted endmembers with the reference spectra from the USGS library [12]. The problem that arises immediately is the correct choice of reference material. The original papers [10] and [11] from which Fig. 5.5 is taken are not very clear on which exact reference sample is used for matching. For example there are 17 different Kaolinite mineral samples and 23 Kaolinite mixtures in the library, and choosing the corresponding “Kaolinite + smectite or muscovite” reference sample is not evident. The reference samples we have chosen as reference endmembers are all collected at the Cuprite scene, have

B or higher purity level (i.e. have slight impurities) and are analyzed with the 2151-band ASD Fieldspec spectrometer. These are CU-98-5C (Alunite), CU91-252D (Muscovite) and CU91-200A (Kaolinite). These reference endmember spectra are depicted in Fig. 5.7d.

Once we have selected the reference endmembers for performance tests, we also need the reference abundance maps. For this we use the modified ICE algorithm with $\mu = 0.01$ and \mathbf{E} fixed to the aforementioned USGS endmembers. We introduce a new parameter $\mathbf{a} \in \mathbb{R}^M$ which we interpret as the gain factor for each endmember. This is needed since spectra of materials measured remotely are almost always much weaker than those of pure reference materials. Minimizing the objective function of the ICE algorithm results in $\widehat{\mathbf{W}}$ and $\widehat{\mathbf{a}}$. The reference abundance maps are depicted in Fig. 5.7a-c.

5.4 Results

The resulting endmembers and abundance maps are shown in Fig. 5.7. A summary of the SAM errors is given in Table 5.6. Do note that these error measures are here for the sake of completeness and do not reflect performance correctly.²

Both ICE and ICE-S perform similarly here given the same volume regularization. For high volume regularization, the algorithm assumes more pure endmembers, which results in higher contrast abundance maps. The unregularized BayesNMF-Vol seems to perform badly. Volume regularization seems to have a large effect here – in contrast to our tests with the synthetic image.

	Parameters			Error metric
	μ	ν	γ	SAM
ICE	0.01	-	-	0.41
ICE-S	0.5	0.02	-	0.43
BayesNMF-Vol	-	-	0	0.92
ICE	0.005	-	-	0.36
ICE-S	0.0015	0.55	-	0.37
BayesNMF-Vol	-	-	0.001	0.43

Figure 5.6: SAM comparison for different parameter configurations.

²First, there is the uncertainty about the real endmembers in the image. As already noted, the spectra of reference materials is taken from small samples in labs within a controlled environment, while we are using spectra of some 18m² of rock from an altitude of 20 km.

Second, the identification of materials based on (endmember) spectra is a topic on its own. Clark [11] notes that trusting in the least-squares matching algorithms as a form of similarity between reference and unknown is far from sufficient. Materials can be spectrally similar but chemically very different. Although not completely identical, they are similar enough “that noise and natural variations in field spectra make the assignment of the proper threshold at which to define identification or misidentification problematic.” So simple measures like MSE (and even SAM) might be deceiving.

In short, assessing the accuracy of the estimated endmembers would require an experienced geologist/spectroscopist with a specific methodology. USGS Tetracorder [11] is an expert system that tries to mimic that methodology. Unfortunately, explaining and assessing the estimated endmembers with such a methodology would require a thesis on its own.

5.5 Conclusion

In this chapter we analyzed the “alunite hill” that appears in the AVIRIS Cuprite scene. For reference endmembers we based ourselves on geological work in the area and the USGS library. The reference abundance maps we extracted using a modified ICE algorithm based on the selected reference spectra. The SAM error results rate the ICE algorithm with $\mu = 0.005$ as the best performer.

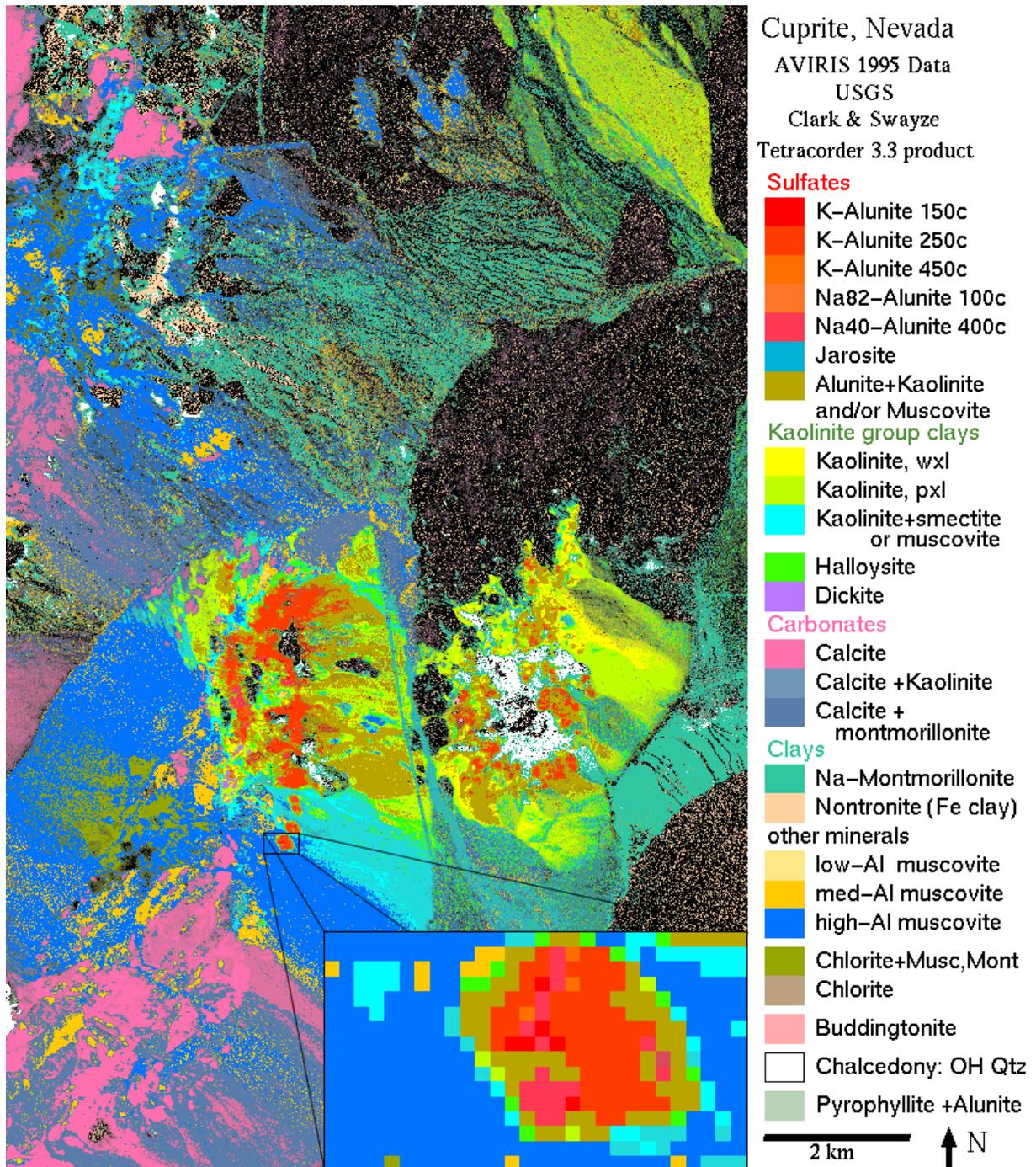


Figure 5.5: The image is taken from [11] with the “alunite hill” (16×28 pixels) zoomed in at the bottom right for clarity. The materials are identified using the USGS spectral library and Tetracoder’s least squares shape-matching algorithm. According to [11] “mineral maps such as these have been extensively field checked to confirm the accuracy of the algorithm”. Note that the surroundings (blue) are mostly high-Al muscovite. The hill edges are mainly kaolinite mixed with muscovite (cyan) and/or alunite (olive green), to mainly K-Alunite (250C - medium temperature) towards the center (orange).

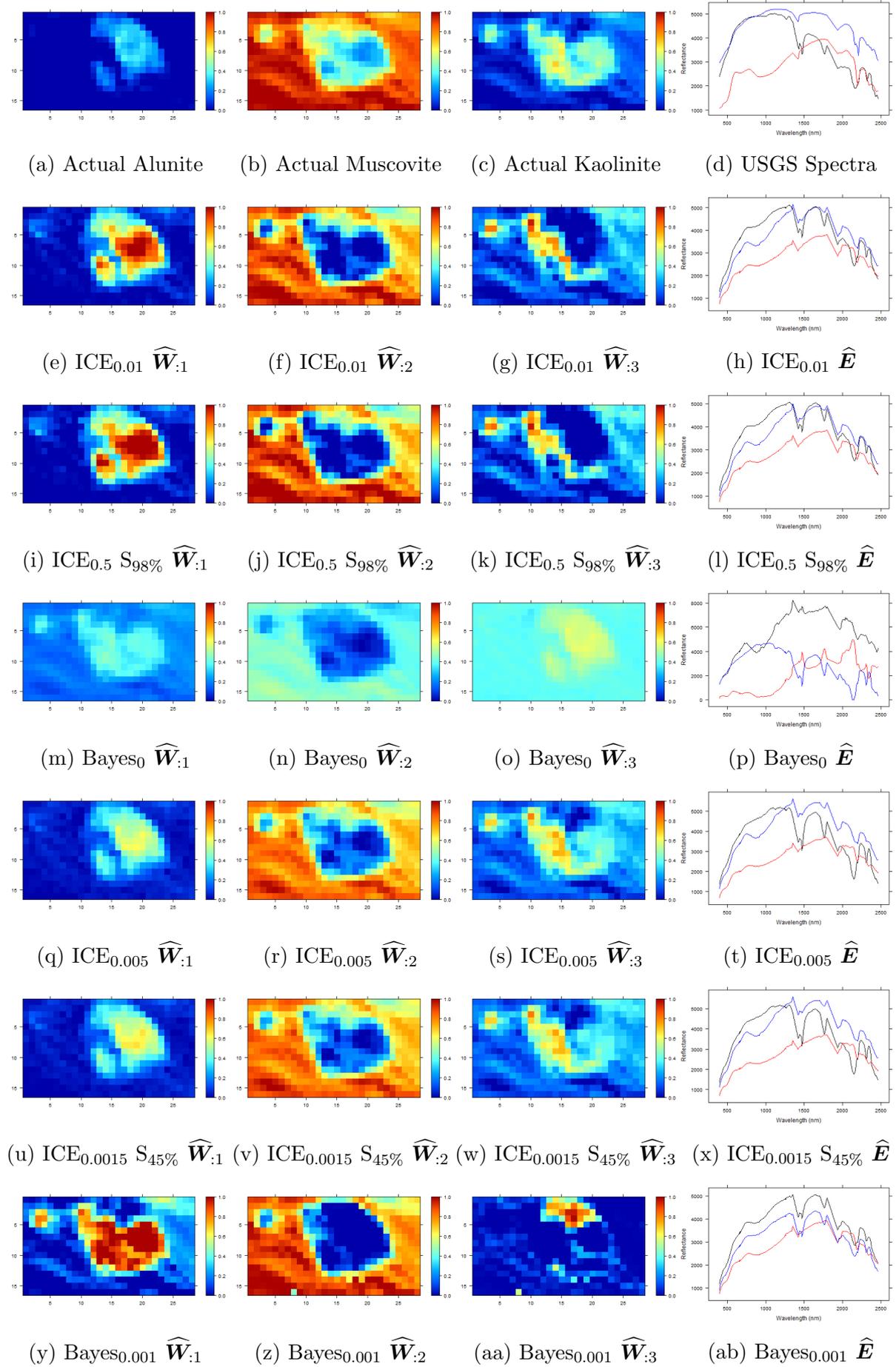


Figure 5.7: The algorithms and parameters used here are the same as in Fig. 4.2. The endmembers are Alunite (black), Muscovite (red) and Kaolinite (blue).

Chapter 6

Conclusion

In this thesis we explained how hyperspectral images are made and their main characteristic, the large amount of spectral bands. We defined a linear mixing model (LMM) which was driven by the idea that due to low spatial resolution of a hyperspectral sensor, spectra of many materials are mixed within each pixel. The LMM allows us to do the inverse step in a mathematically sound way: extract the material spectra from each pixel spectrum. For a whole hyperspectral image, this results in a few constituting materials – of which the spectra are called endmembers – and abundance maps of these materials. The process is called hyperspectral unmixing.

We briefly discussed three kinds of approaches to hyperspectral unmixing: geometrical, statistical and sparse. In this thesis we focused on Berman’s ICE algorithm which is typifying for the geometrical approach. We also focused on one of the first Bayesian approaches, called Arngren’s BayesNMF-Vol algorithm. Since the Bayesian approaches are currently a hot topic, we also sketched the complexity of some state-of-the-art approaches compared to the one described in this thesis.

Our first focus was the ICE algorithm. We first discussed the objective function derivation and implementation details since the original papers are quite succinct on that part. This resulted in small improvements such as the closed form solution for \mathbf{E}^* and a stability improvement of the objective function by dividing it with the number of bands B . Subsequently we added spatial regularization to the ICE algorithm. We dubbed this the ICE-S algorithm. Spatial regularization favors structure in abundance maps, while disfavoring randomness. The synthetic and real data experiments showed clearly the smoothing effect of ICE-S, but were not able to show any difference in endmember estimates – compared to the ICE algorithm under the same volume regularization.

Our second focus was the BayesNMF-Vol algorithm. In a similar way as for the ICE algorithm, we analyzed the assumptions and implementation details in much greater detail than the original papers. We showed how the assumption of a DAG is essential for deriving the conditional distributions used by the Gibbs sampler. We also showed in full detail how sampling from a degenerate distribution can be done – something that proved quite challenging but was only mentioned as a side note in the original papers. We also explained the relationship between the Bayesian MAP estimators and ICE estimators. This led us to believe that the two algorithms are similar, and that the Bayesian approach should also be volume-regularized. The latter was proven wrong for at least two reasons. For one, due to high dimensionality one does not focus on the joint MAP estimator, but on the marginal MAP estimators. Second, because the Bayesian approach is intrinsically volume-regularized due to the “uninformative” prior for the abundances.

The ICE, ICE-S and BayesNMF-Vol algorithms were tested on both a synthetic and a real image. ICE-S was able to produce the best endmembers on the synthetic data once the hyper-parameters were optimized. This was expected because of the extra parameter (and thus an extra degree of freedom). But the values were very dependent on the initialization values. So much in fact that the outcomes were rather unreliable. We visualized this uncertainty by sampling SAM values over the whole hyper-parameter space. Overall we feel that the ICE and ICE-S perform similarly under the same volume regularization. The spatial regularization can be used in such cases to smoothen the abundance maps. On real data with Berman's suggested parametrization the ICE algorithm proved best. What we noticed almost consistently during our tests is the bad performance of the BayesNMF-Vol algorithm when not using volume regularization. A small volume regularization often helps the algorithm to achieve (much) better results, especially in higher dimensions.

This thesis is accompanied by a framework in R for hyperspectral analysis based on the *hyperSpec* package. Together with the implementation of the three algorithms, various functions were made for random image generation, false-color-, abundance- and endmember-plotting, parallel hyper-parameter optimization and USGS library spectra extraction – among others – to streamline the analysis process. Lots remains to be done. The research of hyperspectral unmixing is a never ending endeavor, mainly because of underspecification. But the framework can be used as a good starting point for anyone who wishes to start implementing new unmixing algorithms in R.

Appendices

Appendix A

ICE appendix

A.1 Relation between V and SSD

This is used in Eq. 2.9. Berman [5] notes these equivalent relations, but never derives them. We start from the sum of variance of \mathbf{e}_b and show that it is proportional to the sum of squared distances (SSD) or rewrite it in matrix form which can be used in Eq. 2.13 to solve for \mathbf{E}^* .

$$\begin{aligned}
 V &= \sum_{b=1}^B Var(\mathbf{e}_b) \\
 &= \sum_{b=1}^B \frac{\sum_{m=1}^M \left[e_{bm} - \frac{\sum_{m=1}^M e_{bm}}{M} \right]^2}{M-1} \\
 &= \frac{\sum_{b=1}^B \sum_{m=1}^M e_{bm}^2}{M-1} - \frac{2 \sum_{b=1}^B \sum_{m=1}^M \left(e_{bm} \sum_{m=1}^M e_{bm} \right)}{M(M-1)} + \frac{\sum_{b=1}^B \sum_{m=1}^M \left(\sum_{m=1}^M e_{bm} \right)^2}{M^2(M-1)} \\
 &= \frac{\sum_{b=1}^B \sum_{m=1}^M e_{bm}^2}{M-1} - \frac{2 \sum_{b=1}^B \left(\sum_{m=1}^M e_{bm} \right)^2}{M(M-1)} + \frac{\sum_{b=1}^B \left(\sum_{m=1}^M e_{bm} \right)^2}{M(M-1)} \\
 &= \frac{\sum_{b=1}^B \sum_{m=1}^M e_{bm}^2}{M-1} - \frac{\sum_{b=1}^B \left(\sum_{m=1}^M e_{bm} \right)^2}{M(M-1)}
 \end{aligned} \tag{A.1}$$

From here we can either go towards SSD:

$$\begin{aligned}
V &= \frac{M \sum_{m=1}^M \mathbf{e}_m^T \mathbf{e}_m - \left(\sum_{m=1}^M \mathbf{e}_m \right)^T \left(\sum_{m=1}^M \mathbf{e}_m \right)}{M(M-1)} \\
&= \frac{\sum_{m=1}^M \sum_{k=1}^M (\mathbf{e}_k^T \mathbf{e}_k + \mathbf{e}_m \mathbf{e}_m - 2\mathbf{e}_k \mathbf{e}_m) / 2}{M(M-1)} \\
&= \frac{\sum_{m=1}^M \sum_{k=1}^M (\mathbf{e}_k - \mathbf{e}_m)^T (\mathbf{e}_k - \mathbf{e}_m) / 2}{M(M-1)} \\
&= \frac{\sum_{m=1}^{M-1} \sum_{k=m+1}^M (\mathbf{e}_k - \mathbf{e}_m)^T (\mathbf{e}_k - \mathbf{e}_m)}{M(M-1)} \\
&= \frac{SSD}{M(M-1)}
\end{aligned} \tag{A.2}$$

Or we can rewrite it in function of \mathbf{e}_b :

$$\begin{aligned}
V &= \frac{\sum_{b=1}^B [M \mathbf{e}_b^T \mathbf{e}_b - \mathbf{e}_b^T \mathbf{1} \mathbf{1}^T \mathbf{e}_b]}{M(M-1)} \\
&= \sum_{b=1}^B \frac{\mathbf{e}_b^T (\mathbf{I}_M - \frac{\mathbf{1} \mathbf{1}^T}{M}) \mathbf{e}_b}{M-1}
\end{aligned} \tag{A.3}$$

A.2 Minimum of $L_{reg}(\mathbf{E})$ given \mathbf{W}

This is the solution for the minimization of Eq. 2.3.4 using differential calculus [21]. We start by calculating the differential. Then we massage it into canonical form $dL_{reg}(\mathbf{E}) = \text{tr}(\mathbf{A}^T d\mathbf{E})$ from which we deduce that $DL_{reg}(\mathbf{E}) = \mathbf{A}^T$. Once we know the derivative we can set it equal to 0 and solve for \mathbf{E} .

$$\begin{aligned}
dL_{reg}(\mathbf{E}) &= d \frac{(1-\mu)}{N} \text{tr} [(\mathbf{X} - \mathbf{W} \mathbf{E}^T)^T (\mathbf{X} - \mathbf{W} \mathbf{E}^T)] \\
&\quad + d \frac{\mu}{M-1} \text{tr} [\mathbf{E} (\mathbf{I}_M - \mathbf{1} \mathbf{1}^T / M) \mathbf{E}^T] \\
&= \frac{(1-\mu)}{N} \text{tr} [-2\mathbf{X}^T \mathbf{W} d\mathbf{E}^T + 2\mathbf{E} \mathbf{W}^T \mathbf{W} d\mathbf{E}^T] \\
&\quad + \frac{\mu}{M-1} \text{tr} [2\mathbf{E} d\mathbf{E}^T - 2\mathbf{E} \mathbf{1} \mathbf{1}^T / M d\mathbf{E}^T] \\
&= \text{tr} \left[\frac{2(1-\mu)}{N} (-\mathbf{X}^T \mathbf{W} + \mathbf{E} \mathbf{W}^T \mathbf{W} + \lambda(\mathbf{E} - \mathbf{E} \mathbf{1} \mathbf{1}^T / M)) d\mathbf{E}^T \right] \\
&= \text{tr} \left[\frac{2(1-\mu)}{N} (-\mathbf{W}^T \mathbf{X} + \mathbf{W}^T \mathbf{W} \mathbf{E}^T + \lambda(\mathbf{I}_M - \mathbf{1} \mathbf{1}^T / M) \mathbf{E}^T) d\mathbf{E} \right] \\
&= \text{tr} [\mathbf{A}^T d\mathbf{E}]
\end{aligned} \tag{A.4}$$

where $\lambda = N\mu / [(1-\mu)(M-1)]$. We reduced $dL_{reg}(\mathbf{E})$ to the canonical form from which we deduce that $DL_{reg}(\mathbf{E}) = \mathbf{A}^T$. Now we can solve $\mathbf{A}^T = 0$ for \mathbf{E}^* :

$$\arg \min_{\mathbf{E}} L_{reg}(\mathbf{E}) = \mathbf{X}^T \mathbf{W} [\mathbf{W}^T \mathbf{W} + \lambda(\mathbf{I}_M - \mathbf{1} \mathbf{1}^T / M)]^{-1} \tag{A.5}$$

A.3 Derivation of $S_n(\mathbf{w}_n)$

We use here the definitions from Section 2.5. We further define l as the rows of $\mathbf{U}^{[n]}$ corresponding to the elements in \mathcal{K}_n .

$$\begin{aligned}
S_n(\mathbf{w}_n) &= \sum_{k \in \mathcal{K}} S_k(\mathbf{w}_n) + C_1^{te} \\
&= \sum_k \sum_m \text{Var}(\mathbf{U}_{:m}^{[k]}) + C_1^{te} \\
&= \sum_k \sum_m \sum_l \frac{(w_{lm} - \sum_l w_{lm}/K_k)^2}{K_k - 1} + C_1^{te} \\
&= \sum_k \sum_m \left[\frac{\sum_l w_{lm}^2}{K_k - 1} - \frac{(\sum_l w_{lm})^2}{K_k(K_k - 1)} \right] + C_1^{te} \\
&= \sum_k \sum_m \left[\frac{w_{nm}^2 + \sum_{l \neq n} w_{lm}^2}{K_k - 1} - \frac{w_{nm}^2 + 2w_{nm} \sum_{l \neq n} w_{lm} + (\sum_{l \neq n} w_{lm})^2}{K_k(K_k - 1)} \right] + C_1^{te} \\
&= \sum_k \sum_m \left[\frac{w_{nm}^2}{K_k - 1} - \frac{w_{nm}^2 + 2w_{nm} \sum_{l \neq n} w_{lm}}{K_k(K_k - 1)} \right] + C_2^{te} \\
&= \sum_{k \in \mathcal{K}} \frac{\mathbf{w}_n^T \mathbf{w}_n}{K_k} - 2\mathbf{w}_n^T \sum_{k \in \mathcal{K}} \frac{(\mathbf{U}_{n,:}^k)^T \mathbf{1}}{K_k(K_k - 1)} + C_2^{te}
\end{aligned} \tag{A.6}$$

Appendix B

BayesNMF-Vol appendix

B.1 Density of linearly transformed variable

This is the deduction of the not so trivial step in Section 3.2.4. We have the linear transformation of $\sigma' = \phi(\sigma) = c\sigma$. We further assume that $\sigma \sim f(\sigma)$. We are looking for $f'(\sigma')$. We start by writing the probability

$$\begin{aligned} P(\sigma' < a) &= \int_0^a f'(\sigma') d\sigma' \\ &= \int_0^{\phi^{-1}(a)} f(\sigma) d\sigma \\ &= \int_0^a f(\phi^{-1}(\sigma')) \frac{d\phi^{-1}(\sigma')}{d\sigma'} d\sigma' \\ &= \int_0^a f(\phi^{-1}(\sigma')) \frac{1}{c} d\sigma' \end{aligned} \tag{B.1}$$

where we use integration by substitution in the second last equation. Thus we have that $f'(\sigma') = f(\phi^{-1}(\sigma'))/c = f(\sigma)/c$.

B.2 Derivation of $\mathcal{N}(x_k | \mathbf{x}_{\tilde{k}}, \hat{\mu}_k, \hat{\sigma}_k^2)$ from $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

This form is mainly useful for the Gibbs sampler where we wish to sample x_k from \mathbf{x} , one at a time. Suppose that $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ and that $\boldsymbol{\Sigma}^{-1} = (\sigma_{ij})$.

$$\begin{aligned}
f(x_k | \mathbf{x}_{\tilde{k}}) &= \frac{\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\int_{x_k} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) dx_k} \\
&= C_1^{te} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \\
&= C_2^{te} \exp \left[-\frac{1}{2} (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right] \\
&= C_2^{te} \exp \left[-\frac{1}{2} \left(\sum_i x_i \sum_j \sigma_{ij} y_j - 2 \sum_i x_i \sum_j \sigma_{ij} \mu_j \right) \right] \\
&= C_3^{te} \exp \left[-\frac{1}{2} \left(\sigma_{kk} x_k^2 + x_k \sum_{j \neq k} \sigma_{kj} y_j + x_k \sum_{i \neq k} \sigma_{ik} y_i - 2x_k \sum_j \sigma_{kj} \mu_j \right) \right] \\
&= C_3^{te} \exp \left[-\frac{1}{2} \left(\sigma_{kk} x_k^2 + x_k \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \mathbf{y}_{\tilde{k}} + x_k (\boldsymbol{\Sigma}_{\tilde{k}:k}^{-1})^T \mathbf{y}_{\tilde{k}} - 2x_k \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \boldsymbol{\mu} \right) \right] \tag{B.2} \\
&= C_3^{te} \exp \left[-\frac{1}{2} \frac{\left(x_k^2 - 2 \{ \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \boldsymbol{\mu} - \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \mathbf{y}_{\tilde{k}} / 2 - (\boldsymbol{\Sigma}_{\tilde{k}:k}^{-1})^T \mathbf{y}_{\tilde{k}} / 2 \} \sigma_{kk}^{-1} x_k \right)}{\sigma_{kk}^{-1}} \right] \\
&= C_4^{te} \exp \left[-\frac{1}{2} \frac{\left(x_k - \left\{ \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \boldsymbol{\mu} - \frac{1}{2} \left[\boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} + (\boldsymbol{\Sigma}_{\tilde{k}:k}^{-1})^T \right] \mathbf{y}_{\tilde{k}} \right\} \sigma_{kk}^{-1} \right)^2}{\sigma_{kk}^{-1}} \right] \\
&= C_4^{te} \exp \left[-\frac{1}{2} \frac{\left(x_k - \left\{ \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \boldsymbol{\mu} - \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \mathbf{y}_{\tilde{k}} \right\} \sigma_{kk}^{-1} \right)^2}{\sigma_{kk}^{-1}} \right] \\
&= \mathcal{N}(x_k | \hat{\mu}_k, \hat{\sigma}_k^2)
\end{aligned}$$

The second to last equality follows from the fact that $\boldsymbol{\Sigma}$ is symmetric. From the above we conclude that

$$\hat{\sigma}_k^2 = (\boldsymbol{\Sigma}_{kk}^{-1})^{-1} \tag{B.3}$$

$$\hat{\mu}_k = \frac{\boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \boldsymbol{\mu} - \boldsymbol{\Sigma}_{k:\tilde{k}}^{-1} \mathbf{y}_{\tilde{k}}}{\boldsymbol{\Sigma}_{kk}^{-1}} \tag{B.4}$$

B.3 Derivation of $f_{W|E, \sigma^2, X}$

These results are used in Section 3.4.4.

$$\begin{aligned}
& f(\mathbf{W}|\mathbf{E}, \sigma^2, \mathbf{X}, \alpha, \beta, \gamma) \\
&= \frac{f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2|\alpha, \beta, \gamma)}{\int_{\mathbf{W}} f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2|\alpha, \beta, \gamma) d\mathbf{W}} \\
&= \frac{f(\mathbf{X}|\mathbf{E}, \mathbf{W}, \sigma^2)f(\mathbf{W})}{\int_{\mathbf{W}} f(\mathbf{X}|\mathbf{E}, \mathbf{W}, \sigma^2)f(\mathbf{W}) d\mathbf{W}} \\
&= f(\mathbf{W}|\mathbf{E}, \sigma^2, \mathbf{X}) \\
&\propto \prod_{n=1}^N \left[\mathcal{N}(\mathbf{x}_n|\mathbf{E}\mathbf{w}_n, \sigma^2) \mathbb{I}[\|\mathbf{w}_n\|_1 = 1] \prod_{m=1}^M \mathbb{I}[w_{nm} \geq 0] \right] \\
&\propto \prod_{n=1}^N \left[\mathcal{N}(\mathbf{w}_n|\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_w) \mathbb{I}[\|\mathbf{w}_n\|_1 = 1] \prod_{m=1}^M \mathbb{I}[w_{nm} \geq 0] \right] \\
&\propto \prod_{n=1}^N f(\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n)
\end{aligned} \tag{B.5}$$

Normal density function $f_{\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n}$ from Eq. B.5 must be rewritten in function of \mathbf{w}_n . Thus we equate

$$\begin{aligned}
\mathcal{N}(\mathbf{x}_n|\mathbf{E}\mathbf{w}_n, \sigma^2) &= \frac{\exp\left[-\frac{1}{2}(\mathbf{x}_n - \mathbf{E}\mathbf{w}_n)^T(\sigma^2\mathbf{I})^{-1}(\mathbf{x}_n - \mathbf{E}\mathbf{w}_n)\right]}{\sqrt{|2\pi\sigma^2\mathbf{I}|}} \\
&\propto \frac{\exp\left[-\frac{1}{2}(\mathbf{w}_n - \boldsymbol{\mu}_n)^T\boldsymbol{\Sigma}_w^{-1}(\mathbf{w}_n - \boldsymbol{\mu}_n)\right]}{\sqrt{|2\pi\boldsymbol{\Sigma}_w|}} \\
&= \mathcal{N}(\mathbf{w}_n|\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_w)
\end{aligned} \tag{B.6}$$

from which we get parameters $\boldsymbol{\mu}_n$ and $\boldsymbol{\Sigma}$.

$$\boldsymbol{\Sigma}_w^{-1} = \frac{\mathbf{E}^T \mathbf{E}}{\sigma^2} \tag{B.7}$$

$$\boldsymbol{\mu}_n = (\mathbf{E}^T \mathbf{E})^\dagger \mathbf{E}^T \mathbf{x}_n \tag{B.8}$$

The same derivation approach is used in [2].

B.4 Mapping $f_{\mathbf{w}_n|\mathbf{E}, \sigma^2, \mathbf{x}_n}$ to $f_{\mathbf{y}|\mathbf{E}, \sigma^2, \mathbf{x}_n}$

This is the mapping from the de \mathbf{w}_n -space to the \mathbf{y} -subspace in which the equality constraint $\|\mathbf{w}_n\|_1 = 1$ holds, as initiated in Section. 3.4.4. We first introduce a transformed variable \mathbf{y} which relates to \mathbf{w}_n in the following way:

$$\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ -1 & -1 & -1 & \cdots & -1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{M-1} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} w_{n1} \\ w_{n1} \\ \vdots \\ w_{nM-1} \\ w_{nM} \end{bmatrix} \tag{B.9}$$

We write this in a more concise way as:

$$\mathbf{A}\mathbf{y} + \mathbf{b} = \mathbf{w}_n \quad (\text{B.10})$$

Now we can use this equality to derive a distribution of \mathbf{y} based on $f_{\mathbf{w}_n|\mathbf{E},\sigma^2,\mathbf{x}_n}$:

$$\begin{aligned} f(\mathbf{y}|\mathbf{E}, \sigma^2, \mathbf{x}_n) &= f(\mathbf{A}\mathbf{y} + \mathbf{b}|\mathbf{E}, \sigma^2, \mathbf{x}_n) \\ &= \mathcal{N}(\mathbf{A}\mathbf{y} + \mathbf{b}|\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_w) \mathbb{I}[\|\mathbf{A}\mathbf{y} + \mathbf{b}\|_1 = 1] \prod_{m=1}^M \mathbb{I}[(\mathbf{A}\mathbf{y} + \mathbf{b})_m \geq 0] \\ &= \mathcal{N}(\mathbf{A}\mathbf{y} + \mathbf{b}|\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_w) \prod_{m=1}^M \mathbb{I}[(\mathbf{A}\mathbf{y} + \mathbf{b})_m \geq 0] \\ &= \mathcal{N}(\mathbf{A}\mathbf{y} + \mathbf{b}|\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_w) \mathbb{I}[\|\mathbf{y}\|_1 \leq 1] \prod_{m=1}^{M-1} \mathbb{I}[y_m \geq 0] \\ &= \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \mathbb{I}[\|\mathbf{y}\|_1 \leq 1] \prod_{m=1}^{M-1} \mathbb{I}[y_m \geq 0] \end{aligned} \quad (\text{B.11})$$

where

$$\boldsymbol{\Sigma}_y^{-1} = \mathbf{A}^T \boldsymbol{\Sigma}_w^{-1} \mathbf{A} = \frac{\mathbf{A}^T \mathbf{E}^T \mathbf{E} \mathbf{A}}{\sigma^2} \quad (\text{B.12})$$

$$\begin{aligned} \boldsymbol{\mu}_y &= \boldsymbol{\Sigma}_y \mathbf{A}^T \boldsymbol{\Sigma}_w^{-1} (\boldsymbol{\mu}_n - \mathbf{b}) \\ &= (\mathbf{A}^T \mathbf{E}^T \mathbf{E} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{E}^T \mathbf{E} [(\mathbf{E}^T \mathbf{E})^+ \mathbf{E}^T \mathbf{x}_n - \mathbf{b}] \\ &= (\mathbf{A}^T \mathbf{E}^T \mathbf{E} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{E}^T [\mathbf{x}_n - \mathbf{E}\mathbf{b}] \end{aligned} \quad (\text{B.13})$$

To get rid of the MP inverse we use the equality $(\mathbf{A}^T \mathbf{A})^\dagger \mathbf{A}^T = \mathbf{A}^\dagger$ and $\mathbf{A}^T \mathbf{A} \mathbf{A}^\dagger = \mathbf{A}^T$ which holds for any matrix \mathbf{A} (cf. [21, p.38]).

The $(M-1, M-1)$ matrix $\mathbf{A}^T \mathbf{E}^T \mathbf{E} \mathbf{A}$ is nonsingular \iff columns of $\mathbf{E}\mathbf{A}$ are independent $\iff \text{rank}(\mathbf{E}\mathbf{A}) = M-1 \iff \ker(\mathbf{E}\mathbf{A}) = \{\mathbf{x} : \mathbf{A}\mathbf{x} \in \ker \mathbf{E}\} = \{\mathbf{0}\}$. Knowing that $\text{rank}(\mathbf{E}\mathbf{A}) \leq \min(\text{rank}(\mathbf{E}), M-1) \leq \min(B, M-1)$ the matrix will certainly be singular if $B < M-1$. More generally, it will be singular if $\text{range}(\mathbf{A})$ intersects with $\ker(\mathbf{E})$ in more than $\{\mathbf{0}\}$. From our testing, the latter seems to be very unlikely, thus the matrix can be assumed nonsingular.

B.5 Derivation of $f_{y_k|\mathbf{y}_{\setminus k}, \mathbf{E}, \sigma^2, \mathbf{x}_n}$ for Gibbs sampling

We start from results given in Eq. 3.18.

$$\begin{aligned}
& f(y_k | \mathbf{y}_{\tilde{k}}, \mathbf{E}, \sigma^2, \mathbf{x}_n) \\
&= \frac{f(\mathbf{y} | \mathbf{E}, \sigma^2, \mathbf{x}_n)}{\int_{y_k} f(\mathbf{y} | \mathbf{E}, \sigma^2, \mathbf{x}_n) dy_k} \\
&= \frac{\mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_{\mathbf{y}}, \boldsymbol{\Sigma}_{\mathbf{y}}) \mathbb{I}[\|\mathbf{y}\|_1 \leq 1] \prod_{m=1}^{M-1} \mathbb{I}[y_m \geq 0]}{\prod_{m \neq k}^{M-1} \mathbb{I}[y_m \geq 0] \int_{y_k} \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_{\mathbf{y}}, \boldsymbol{\Sigma}_{\mathbf{y}}) \mathbb{I}[\|\mathbf{y}\|_1 \leq 1] \mathbb{I}[y_k \geq 0] dy_k} \\
&= \frac{\mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_{\mathbf{y}}, \boldsymbol{\Sigma}_{\mathbf{y}}) \mathbb{I}\left[y_k \leq 1 - \sum_{m \neq k}^{M-1} y_m\right] \mathbb{I}[y_k \geq 0]}{\int_{y_k} \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_{\mathbf{y}}, \boldsymbol{\Sigma}_{\mathbf{y}}) \mathbb{I}[\|\mathbf{y}\|_1 \leq 1] \mathbb{I}[y_k \geq 0] dy_k} \\
&= \mathcal{N}(y_k | \mu_{yk}, \sigma_{yk}^2) \mathbb{I}\left[y_k \leq 1 - \sum_{m \neq k}^{M-1} y_m\right] \mathbb{I}[y_k \geq 0]
\end{aligned} \tag{B.14}$$

where the parameters μ_k and σ_k^2 we compute from the result in Appendix B.2:

$$\sigma_{yk}^2 = \left[(\boldsymbol{\Sigma}_{\mathbf{y}}^{-1})_{kk} \right]^{-1} \tag{B.15}$$

$$\mu_{yk} = \frac{(\boldsymbol{\Sigma}_{\mathbf{y}}^{-1})_{k:} \boldsymbol{\mu}_{\mathbf{y}} - (\boldsymbol{\Sigma}_{\mathbf{y}}^{-1})_{k:\tilde{k}} \mathbf{y}_{\tilde{k}}}{\sigma_{yk}^{-2}} \tag{B.16}$$

Note that we are generating a Gibbs sequence for \mathbf{y} , so we need previous components of \mathbf{y} . We have \mathbf{w}_n though, and we can get \mathbf{y} by the inverse transformation of Eq. B.10. Matrix \mathbf{A} is a $(M \times M - 1)$ matrix. It has no inverse since it is non-square. But it does have a left inverse \mathbf{A}_{Left}^{-1} since its columns are independent and equivalently that the Gram matrix $\mathbf{A}^T \mathbf{A}$ is non-singular. We use this property to write \mathbf{y} in function of \mathbf{w}_n :

$$\mathbf{y} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{A} \mathbf{y} = \mathbf{A}_{Left}^{-1} \mathbf{A} \mathbf{y} = \mathbf{A}_{Left}^{-1} (\mathbf{w}_n - \mathbf{b}) = \mathbf{w}_{n\tilde{M}} \tag{B.17}$$

Note the last equality. \mathbf{y} is just \mathbf{w}_n with the last component removed!

B.6 Derivation of $f_{E|W, \sigma^2, \mathbf{X}, \gamma}$

We start from the results given in Eq. 3.21.

$$\begin{aligned}
& f(\mathbf{E}|\mathbf{W}, \sigma^2, \mathbf{X}, \alpha, \beta, \gamma) \\
&= \frac{f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2|\alpha, \beta, \gamma)}{\int_{\mathbf{E}} f_{\text{BayesNMF-Vol}}(\mathbf{X}, \mathbf{E}, \mathbf{W}, \sigma^2|\alpha, \beta, \gamma) d\mathbf{E}} \\
&= \frac{f(\mathbf{X}|\mathbf{E}, \mathbf{W}, \sigma^2)f(\mathbf{E}|\gamma)}{\int_{\mathbf{E}} f(\mathbf{X}|\mathbf{E}, \mathbf{W}, \sigma^2)f(\mathbf{E}|\gamma) d\mathbf{E}} \\
&= f(\mathbf{E}|\mathbf{W}, \sigma^2, \mathbf{X}, \gamma) \\
&\propto \prod_{n=1}^N \prod_{b=1}^B \frac{\exp\left[-\frac{1}{2\sigma^2}(x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2\right]}{\sqrt{2\pi\sigma^2}} e^{-\gamma V(\mathbf{E})} \prod_{b=1}^B \prod_{m=1}^M \mathbb{I}[e_{bm} \geq 0] \\
&= \prod_{n=1}^N \prod_{b=1}^B \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_{nb} - \mathbf{w}_n^T \mathbf{e}_b)^2\right) \\
&\quad \times \exp\left(-\gamma \sum_{b=1}^B \frac{\mathbf{e}_b^T (\mathbf{I}_M - \frac{\mathbf{1}\mathbf{1}^T}{M}) \mathbf{e}_b}{M-1}\right) \prod_{b=1}^B \prod_{m=1}^M \mathbb{I}[e_{bm} \geq 0] \\
&= \prod_{b=1}^B \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^N \exp\left(-\frac{1}{2} \left[\frac{(\mathbf{x}_b - \mathbf{W}\mathbf{e}_b)^T (\mathbf{x}_b - \mathbf{W}\mathbf{e}_b)}{\sigma^2} + 2\gamma \frac{\mathbf{e}_b^T (\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M) \mathbf{e}_b}{M-1} \right]\right) \\
&\quad \times \prod_{b=1}^B \prod_{m=1}^M \mathbb{I}[e_{bm} \geq 0] \\
&\propto \prod_{b=1}^B \left[\mathcal{N}(\mathbf{e}_b|\boldsymbol{\mu}_b, \boldsymbol{\Sigma}_{\mathbf{E}}) \prod_{m=1}^M \mathbb{I}[e_{bm} \geq 0] \right] \\
&\propto \prod_{b=1}^B f(\mathbf{e}_b|\mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma)
\end{aligned} \tag{B.18}$$

So we need to equate

$$\begin{aligned}
& (\mathbf{e}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Sigma}_{\mathbf{E}}^{-1} (\mathbf{e}_b - \boldsymbol{\mu}_b) + C^{te} \\
&= \frac{(\mathbf{x}_b - \mathbf{W}\mathbf{e}_b)^T (\mathbf{x}_b - \mathbf{W}\mathbf{e}_b)}{\sigma^2} + 2\gamma \frac{\mathbf{e}_b^T (\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M) \mathbf{e}_b}{M-1}
\end{aligned} \tag{B.19}$$

and solve for \mathbf{e}_b and $\boldsymbol{\Sigma}_{\mathbf{E}}^{-1}$ where we get:

$$\boldsymbol{\Sigma}_{\mathbf{E}}^{-1} = \frac{\mathbf{W}^T \mathbf{W}}{\sigma^2} + 2\gamma \left(\frac{\mathbf{I}_M - \mathbf{1}\mathbf{1}^T/M}{M-1} \right) \tag{B.20}$$

$$\boldsymbol{\mu}_b = \frac{\boldsymbol{\Sigma}_{\mathbf{E}} \mathbf{W}^T \mathbf{x}_b}{\sigma^2} \tag{B.21}$$

B.7 Derivation of $f_{e_{bm}|e_{b\tilde{m}}, \mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma}$ for Gibbs sampling

The density functions $f_{e_b|\mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma}$ from Eq. B.18 is relatively difficult to sample from. So we derive $f_{e_{bm}|e_{b\tilde{m}}, \mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma}$ which can be used in a Gibbs sampling procedure.

$$\begin{aligned}
& f(e_{bm} | \mathbf{e}_{b\tilde{m}}, \mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma) \\
&= \frac{f(\mathbf{e}_b | \mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma)}{\int_{e_{bm}} f(\mathbf{e}_b | \mathbf{W}, \sigma^2, \mathbf{x}_b, \gamma) de_{bm}} \\
&= \frac{\mathcal{N}(\mathbf{e}_b | \boldsymbol{\mu}_b, \boldsymbol{\Sigma}_E) \prod_{i=1}^M \mathbb{I}[e_{bi} \geq 0]}{\prod_{i \neq m}^M \mathbb{I}[e_{bi} \geq 0] \int_{e_{bm}} \mathcal{N}(\mathbf{e}_b | \boldsymbol{\mu}_b, \boldsymbol{\Sigma}_E) \mathbb{I}[e_{bm} \geq 0] de_{bm}} \tag{B.22} \\
&= \frac{\mathcal{N}(\mathbf{e}_b | \boldsymbol{\mu}_b, \boldsymbol{\Sigma}_E) \mathbb{I}[e_{bm} \geq 0]}{\int_{e_{bm}} \mathcal{N}(\mathbf{e}_b | \boldsymbol{\mu}_b, \boldsymbol{\Sigma}_E) \mathbb{I}[e_{bm} \geq 0] de_{bm}} \\
&= \mathcal{N}(e_{bm} | \bar{\mu}_{bm}, \bar{\sigma}_m^2) \mathbb{I}[e_{bm} \geq 0]
\end{aligned}$$

Using the results from Appendix B.2 we have

$$\bar{\mu}_{bm} = \frac{(\boldsymbol{\Sigma}_E^{-1})_{m:} \boldsymbol{\mu}_b - (\boldsymbol{\Sigma}_E^{-1})_{m:\tilde{m}} \mathbf{e}_{b\tilde{m}}}{(\boldsymbol{\Sigma}_E^{-1})_{mm}} \tag{B.23}$$

$$\bar{\sigma}_m^2 = 1 / (\boldsymbol{\Sigma}_E^{-1})_{mm} \tag{B.24}$$

Appendix C

Source code

The R source code is in RMD files. The compiled output is available at <http://josipovic.be/statistics/dissertation/>.

All the source code, data, R and its libraries, and output used in this thesis are compiled and made available at: <https://hyperspec-unmix.sourceforge.io>. A small wiki for running the code is available. The source code is composed of RMD files where each covers a topic on its own. For example, the *HyperImage.rmd* file has the definition of the *HyperSpecExt* object, and functions for false color plotting and plotting of the abundance maps. The *Alg-ICE.rmd* has the implementation of the ICE algorithm and depends on, a.o., the definitions in *HyperImage.rmd*. This dependence is declared in the first chunk. This dependence declaration allows *Alg-ICE.rmd* to be compiled on its own. The files with "TEST" suffix, such as *Alg-ICE-TEST.rmd*, contain some kind of processing with the aforementioned algorithm. All RMD files are included in the *Main.rmd* file which, when compiled with *knitr*, gives an overview of all the source code and output.

The R library and used packages – most notably the *HyperSpec* package – are available in the */r* directory. The used Cuprite and USGS data can be found in the */data* directory. These files fall under public license and are made available for convenience and compatibility.

Bibliography

- [1] Morten Arngren, Mikkel N. Schmidt, and Jan Larsen. Bayesian nonnegative matrix factorization with volume prior for unmixing of hyperspectral images. In *2009 IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6. IEEE, 2009.
- [2] Morten Arngren, Mikkel N. Schmidt, and Jan Larsen. Supplementary material for unmixing of hyperspectral images using Bayesian nonnegative matrix factorization with volume prior. 2010.
- [3] Morten Arngren, Mikkel N. Schmidt, and Jan Larsen. Unmixing of hyperspectral images using Bayesian non-negative matrix factorization with volume prior. *Journal of Signal Processing Systems*, 65(3):479–496, 2011.
- [4] David Barber. *Bayesian reasoning and machine learning*. Cambridge University Press, 2012.
- [5] Mark Berman, Harri Kiiveri, Ryan Lagerstrom, Andreas Ernst, Rob Dunne, and Jonathan F. Huntington. Ice: an automated statistical approach to identifying endmembers in hyperspectral images. In *Geoscience and Remote Sensing Symposium, 2003. IGARSS'03. Proceedings. 2003 IEEE International*, volume 1, pages 279–283. IEEE, 2003.
- [6] Mark Berman, Harri Kiiveri, Ryan Lagerstrom, Andreas Ernst, Rob Dunne, and Jonathan F. Huntington. Ice: A statistical approach to identifying endmembers in hyperspectral images. *IEEE transactions on Geoscience and Remote Sensing*, 42(10):2085–2095, 2004.
- [7] Mark Berman, Alope Phatak, Ryan Lagerstrom, and Bayden R. Wood. Ice: a new method for the multivariate curve resolution of hyperspectral images. *Journal of Chemometrics*, 23(2):101–116, 2009.
- [8] José M. Bioucas-Dias, Antonio Plaza, Nicolas Dobigeon, Mario Parente, Qian Du, Paul Gader, and Jocelyn Chanussot. Hyperspectral unmixing overview: Geometrical, statistical, and sparse regression-based approaches. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 5(2):354–379, 2012.
- [9] George Casella and Edward I. George. Explaining the Gibbs sampler. *The American Statistician*, 46(3):167–174, 1992.
- [10] Roger N. Clark, Gregg A. Swayze, and Andrea Gallagher. Mapping minerals with imaging spectroscopy. *US Geological Survey, Office of Mineral Resources Bulletin*, 2039:141–150, 1993.

- [11] Roger N. Clark, Gregg A. Swayze, K. Eric Livo, Raymond F. Kokaly, Steve J. Sutley, J. Brad Dalton, Robert R. McDougal, and Carol A. Gent. Imaging spectroscopy: Earth and planetary remote sensing with the usgs tetracorder and expert systems. *Journal of Geophysical Research: Planets*, 108(E12), 2003.
- [12] Roger N. Clark, Gregg A. Swayze, Richard Wise, K. Eric Livo, T. Hoefen, Raymond F. Kokaly, and Steve J. Sutley. Usgs digital spectral library splib06a. *US Geological Survey, Digital Data Series*, 231, 2007.
- [13] Maurice D. Craig. Minimum-volume transforms for remotely sensed data. *IEEE Transactions on Geoscience and Remote Sensing*, 32(3):542–552, 1994.
- [14] AVIRIS Free Data. Jet propulsion lab. *California Inst. Technol., Pasadena.[Online]*, 1997.
- [15] Philip E. Dennison, Kerry Q. Halligan, and Dar A. Roberts. A comparison of error metrics and constraints for multiple endmember spectral mixture analysis and spectral angle mapper. *Remote Sensing of Environment*, 93(3):359–367, 2004.
- [16] Nicolas Dobigeon, Saïd Moussaoui, Martial Coulon, Jean-Yves Tournet, and Alfred O Hero. Joint Bayesian endmember extraction and linear unmixing for hyperspectral imagery. *IEEE Transactions on Signal Processing*, 57(11):4355–4368, 2009.
- [17] Abderrahim Halimi, Nicolas Dobigeon, Jean-Yves Tournet, and Paul Honeine. A new Bayesian unmixing algorithm for hyperspectral images mitigating endmember variability. In *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*, pages 2469–2473. IEEE, 2015.
- [18] Edwin T. Jaynes. Prior probabilities. *IEEE Transactions on systems science and cybernetics*, 4(3):227–241, 1968.
- [19] Fred A. Kruse, A. B. Lefkoff, J. W. Boardman, K. B. Heidebrecht, A. T. Shapiro, P. J. Barloon, and A. F. H. Goetz. The spectral image processing system (SIPS)-interactive visualization and analysis of imaging spectrometer data. *Remote sensing of environment*, 44(2-3):145–163, 1993.
- [20] Rong Liu, Bo Du, and Liangpei Zhang. Hyperspectral unmixing via double abundance characteristics constraints based nmf. *Remote Sensing*, 8(6):464, 2016.
- [21] Jan R. Magnus and Heinz Neudecker. *Matrix differential calculus with applications in statistics and econometrics*. John Wiley & Sons, 3rd edition, 2007.
- [22] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, New York, 2nd edition, 2006.
- [23] Mikkel Schmidt. Linearly constrained Bayesian matrix factorization for blind source separation. In *Advances in neural information processing systems*, pages 1624–1632, 2009.
- [24] Ben Somers, Gregory P. Asner, Laurent Tits, and Pol Coppin. Endmember variability in spectral mixture analysis: A review. *Remote Sensing of Environment*, 115(7):1603–1616, 2011.

- [25] Wei Tang, Zhenwei Shi, Ying Wu, and Changshui Zhang. Sparse unmixing of hyperspectral data using spectral a priori information. *IEEE Transactions on Geoscience and Remote Sensing*, 53(2):770–783, 2015.
- [26] Michael E. Winter. N-findr: An algorithm for fast autonomous spectral end-member determination in hyperspectral data. In *SPIE's International Symposium on Optical Science, Engineering, and Instrumentation*, pages 266–275. International Society for Optics and Photonics, 1999.
- [27] Stephen J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.



Faculty of Sciences
Department of Applied Mathematics, Computer Science and Statistics
Faculty of Engineering
Department of Telecommunications and Information Processing (TELIN)
Image Processing and Interpretation (IPI)

Volume-based geometric and Bayesian approaches in linear hyperspectral unmixing

Davor JOSIPOVIC

Promoter: Prof. Dr. Aleksandra PIZURICA
Co-promoter: Prof. Dr. Hongyan ZHANG
Co-promoter: Prof. Dr. Dries BENOIT

Thesis submitted in fulfillment of the requirements for the degree of
MASTER OF SCIENCE IN STATISTICAL DATA ANALYSIS

Year: 2016–2017

Compiled: June 21, 2017