

Hierarchical Markovian Models for Joint Classification, Segmentation and Data Reduction of Hyperspectral Images

Nadia BALI, Ali MOHAMMAD-DJAFARI and Adel MOHAMMADPOOR

Laboratoire des Signaux et Systèmes,
Unité mixte de recherche 8506 (CNRS-Supélec-UPS)
Supélec, Plateau de Moulon, 3 rue Joliot Curie, 91192 Gif-sur-Yvette, France.

Abstract. Spectral classification, segmentation and data reduction are the three main problems in hyperspectral image analysis. In this paper we propose a Bayesian estimation approach which tries to give a solution for these three problems jointly. The data reduction problem is modeled as a blind sources separation (BSS) where the data are the m hyperspectral images and the sources are the $n < m$ images which must be mutually the most independent and piecewise homogeneous. To insure these properties, we propose a hierarchical model for the sources with a common hidden classification variable which is modelled via a Potts-Markov field. The joint Bayesian estimation of this hidden variable as well as the sources and the mixing matrix of the BSS problem gives a solution for all the three problems of spectra classification, segmentation and data reduction problems of hyperspectral images. A few simulation results are given to illustrate the performances of the proposed method and some comparison with other classical methods of PCA and ICA used for BSS.

1 Introduction

Hyperspectral images data are often represented either as a set of images $x_\omega(\mathbf{r})$ or as a set of spectra $x_{\mathbf{r}}(\omega)$ where $\omega \in \{1, \dots, \Omega\}$ indexes the wavelength and $\mathbf{r} \in \mathcal{R}$ is a pixel position [1, 2, 3]. In both representations, the data are dependent in both spatial positions and in spectral bands. Classical methods of hyperspectral image analysis try either to classify the spectra $x_\omega(\mathbf{r})$ in K classes $\{s_k(\omega), k = 1, \dots, K\}$ or to classify the images $x_\omega(\mathbf{r})$ in N classes $\{s_j(\mathbf{r}), j = 1, \dots, N\}$, using in both cases, the classical classification methods such as distance based methods (like K -means) or probabilistic methods using the mixture of Gaussian (MoG) modeling of the data. However, these methods either neglect the spatial structure of the spectra or the spectral natures of the pixels along the wavelength bands.

If we consider the data as a set of spectra, then we want to write:

$$x_{\mathbf{r}}(\omega) = \sum_{k=1}^K A_{\mathbf{r},k} s_k(\omega) + \epsilon_{\mathbf{r}}(\omega), \quad (1)$$

where the $s_k(\omega)$ are the K spectral sources and each column of the mixing matrix \mathbf{A} is in fact an image $A_k(\mathbf{r})$. The ideal case here would be to obtain an estimate

for \mathbf{A} such that each column $A_k(\mathbf{r})$ represents an image where only non-zero values for the pixels in the regions which are associated to the spectrum $s_k(\omega)$. At the other hand, if we consider the data as a set of images $x_\omega(\mathbf{r})$, then we have:

$$x_\omega(\mathbf{r}) = \sum_{j=1}^N A_{\omega,j} s_j(\mathbf{r}) + \epsilon_\omega(\mathbf{r}), \quad (2)$$

where the sources $s_j(\mathbf{r})$ are the N source images and each column of the mixing matrix \mathbf{A} in this case correspond to the spectrum $A_j(\omega)$. The ideal case here would be to obtain an estimate for the sources such that the pixels of each image $s_j(\mathbf{r})$ be non-zero only for the positions in the regions which are associated to the spectrum $A_j(\omega)$.

In this paper, we propose to consider the data reduction problem as a blind sources separation (BSS) and use a Bayesian estimation framework with a hierarchical model for the sources with a common hidden classification variable which is modelled via a Potts-Markov field. The joint estimation of this hidden variable as well as the sources and the mixing matrix of the BSS problem gives a solution for all the three problems of spectra classification, segmentation and data reduction problems of hyperspectral images.

2 Proposed data reduction model and method

2.1 Data reduction model

As explained in the introduction, we propose to consider the data reduction problem as in equation (2), written in vector form:

$$\mathbf{x}(\mathbf{r}) = \mathbf{A}\mathbf{s}(\mathbf{r}) + \boldsymbol{\epsilon}(\mathbf{r}) \quad (3)$$

where $\mathbf{x}(\mathbf{r}) = \{x_i(\mathbf{r}), i = 1, \dots, m\}$ is the set of m observed mixed images (hyperspectral images), \mathbf{A} the unknown mixing matrix of dimensions (m, n) , $\mathbf{s}(\mathbf{r}) = \{s_j(\mathbf{r}), j = 1, \dots, n\}$ the set of n unknown components (source images) and $\boldsymbol{\epsilon}(\mathbf{r}) = \{\epsilon_i(\mathbf{r}), i = 1, \dots, m\}$ represents the errors. Now, if we note by $\underline{\mathbf{x}} = \{\mathbf{x}(\mathbf{r}), \mathbf{r} \in \mathcal{R}\}$, $\underline{\mathbf{s}} = \{\mathbf{s}(\mathbf{r}), \mathbf{r} \in \mathcal{R}\}$ and $\underline{\boldsymbol{\epsilon}} = \{\boldsymbol{\epsilon}(\mathbf{r}), \mathbf{r} \in \mathcal{R}\}$, then we can write

$$\underline{\mathbf{x}} = \mathbf{A}\underline{\mathbf{s}} + \underline{\boldsymbol{\epsilon}} \quad (4)$$

In the following, we assume that the errors $\boldsymbol{\epsilon}(\mathbf{r})$ are centered, white, Gaussian with covariance matrix $\boldsymbol{\Sigma}_\epsilon = \text{diag}(\sigma_{\epsilon_1}^2, \dots, \sigma_{\epsilon_m}^2)$. This leads to

$$p(\underline{\mathbf{x}}|\underline{\mathbf{s}}, \mathbf{A}, \boldsymbol{\Sigma}_\epsilon) = \prod_{\mathbf{r}} \mathcal{N}(\mathbf{A}\mathbf{s}(\mathbf{r}), \boldsymbol{\Sigma}_\epsilon) \quad (5)$$

2.2 Sources model

As we mentioned in the introduction, we want to impose to all these sources $\mathbf{s}(\mathbf{r})$ to be piecewise homogeneous and share the same segmentation. This can

be achieved via the introduction of a discrete valued hidden variable $z(\mathbf{r})$ and by assuming the following:

$$p(s_j(\mathbf{r})|z(\mathbf{r}) = k) = \mathcal{N}(m_{j_k}, \sigma_{j_k}^2) \quad (6)$$

and

$$p(z(\mathbf{r}), \mathbf{r} \in \mathcal{R}) \propto \exp \left[\beta \sum_{\mathbf{r} \in \mathcal{R}} \sum_{\mathbf{r}' \in \mathcal{V}(\mathbf{r})} \delta(z(\mathbf{r}) - z(\mathbf{r}')) \right] \quad (7)$$

$z(\mathbf{r})$ then will represent the common segmentation of the sources and the data. We also impose $m_{j_k} = 0$ if $j \neq k$ and $\sigma_{j_k}^2 = .001$ if $j \neq k$ which try to insure that each image $s_j(\mathbf{r})$ be composed of zeros everywhere except those regions associated with class k .

2.3 Data and sources hierarchical model

Combining the data and the sources models of the previous section, we obtain the following hierarchical model:

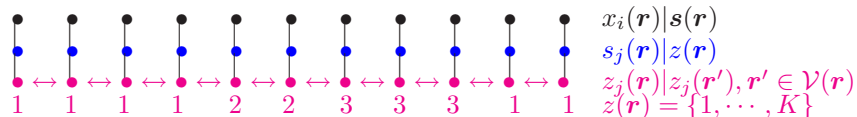


Fig. 1: **Proposed Model:** A hierarchical model for data $x_i(\mathbf{r})$ and sources $s_j(\mathbf{r})$ with hidden common classification and segmentation variables $z(\mathbf{r})$.

3 Bayesian estimation framework

The first step in the Bayesian estimation approach is to find expressions for the likelihood $p(\underline{\mathbf{x}}|\underline{\mathbf{s}}, \mathbf{A}, \mathbf{R}_\epsilon)$ and for the sources as we did in previous section. However, these expressions depend on some hyperparameters $\underline{\boldsymbol{\theta}} = \{\boldsymbol{\theta}_\epsilon, \boldsymbol{\theta}_s\}$ where $\boldsymbol{\theta}_\epsilon = \mathbf{R}_\epsilon$ and $\boldsymbol{\theta}_s = \{(m_{j_k}, \sigma_{j_k}^2)\}$. So, we have to assign $p(\underline{\boldsymbol{\theta}})$ and also $p(\mathbf{A})$. In the following we use conjugate priors for all of them, i.e., Gaussian for the elements of \mathbf{A} , Gaussian for the means m_{j_k} and inverse Gamma for the variances $\sigma_{j_k}^2$ as well as for the noise variances σ_{ϵ_i} .

When all these priors are appropriately assigned, we can obtain an expression for the posterior law

$$p(\underline{\mathbf{s}}, \mathbf{z}, \mathbf{A}, \underline{\boldsymbol{\theta}}|\underline{\mathbf{x}}) \propto p(\underline{\mathbf{x}}|\underline{\mathbf{s}}, \mathbf{A}, \boldsymbol{\theta}_\epsilon) p(\underline{\mathbf{s}}|\mathbf{z}, \boldsymbol{\theta}_s) p(\underline{\boldsymbol{\theta}}) \quad (8)$$

We can then use this posterior law to define an estimator such as Joint Maximum A Posteriori (JMAP) or the Posterior Means (PM). The first needs optimization algorithms and the second integration methods. Both are computationally demanding. Alternate optimization is generally used for the first while the MCMC techniques are used for the second. We propose here the following algorithm:

- First, integrate out $\underline{\mathbf{s}}$ to obtain $p(\mathbf{z}, \mathbf{A}, \underline{\boldsymbol{\theta}}|\underline{\mathbf{x}})$ and then estimate \mathbf{z} by

$$\widehat{\mathbf{z}} \sim p(\mathbf{z}|\widehat{\mathbf{A}}, \widehat{\underline{\boldsymbol{\theta}}}, \underline{\mathbf{x}})$$

- Then, compute the posterior mean of the sources using $p(\underline{\mathbf{s}}|\widehat{\mathbf{z}}, \widehat{\mathbf{A}}, \widehat{\underline{\boldsymbol{\theta}}}, \underline{\mathbf{x}})$.
- Third estimate \mathbf{A} and $\underline{\boldsymbol{\theta}}$ using $p(\mathbf{A}, \underline{\boldsymbol{\theta}}|\underline{\mathbf{s}}, \mathbf{z}, \underline{\mathbf{x}})$

In this algorithm, \sim represents either *argmax* or *generate sample using* or still *compute the Mean Field Approximation (MFA)*. In the following, we give detail expressions of the different conditional laws as well as the details of the proposed algorithm.

4 Expressions of the a posteriori conditional laws

- $p(\underline{\mathbf{s}}|\mathbf{z}, \mathbf{A}, \underline{\boldsymbol{\theta}}, \underline{\mathbf{x}})$:

$$p(\underline{\mathbf{s}}|\mathbf{z}, \underline{\boldsymbol{\theta}}, \underline{\mathbf{x}}) = \prod_{\mathbf{r}} p(\mathbf{s}(\mathbf{r})|\mathbf{z}(\mathbf{r}), \underline{\boldsymbol{\theta}}, \mathbf{x}(\mathbf{r})) = \prod_{\mathbf{r}} \mathcal{N}(\boldsymbol{\mu}(\mathbf{r}), \mathbf{B}(\mathbf{r}))$$

with

$$\begin{cases} \mathbf{B}(\mathbf{r}) = \left[\mathbf{A}^t \boldsymbol{\Sigma}_{\epsilon}^{-1} \mathbf{A} + \boldsymbol{\Sigma}_{z(\mathbf{r})}^{-1} \right]^{-1} \\ \boldsymbol{\mu}(\mathbf{r}) = \mathbf{B}(\mathbf{r}) \left[\mathbf{A}^t \boldsymbol{\Sigma}_{\epsilon}^{-1} \mathbf{x}(\mathbf{r}) + \boldsymbol{\Sigma}_{z(\mathbf{r})}^{-1} \mathbf{m}_{z(\mathbf{r})} \right] \end{cases}$$

$\mathbf{s}(\mathbf{r})$ is then estimated by its *a posteriori* mean so $\bar{\mathbf{s}}(\mathbf{r}) = \boldsymbol{\mu}(\mathbf{r})$

- $p(\mathbf{z}|\mathbf{A}, \underline{\boldsymbol{\theta}}, \underline{\mathbf{x}})$:

$$p(\mathbf{z}|\mathbf{A}, \underline{\boldsymbol{\theta}}, \underline{\mathbf{x}}) \propto p(\underline{\mathbf{x}}|\mathbf{z}, \mathbf{A}, \underline{\boldsymbol{\theta}}) p(\mathbf{z})$$

where

$$p(\underline{\mathbf{x}}|\mathbf{z}, \mathbf{A}, \underline{\boldsymbol{\theta}}) \propto \prod_{\mathbf{r}} p(\mathbf{x}(\mathbf{r})|\mathbf{z}(\mathbf{r}), \mathbf{A}, \underline{\boldsymbol{\theta}}) \propto \prod_{\mathbf{r}} \mathcal{N}(\mathbf{A} \mathbf{m}_{z(\mathbf{r})}, \mathbf{A} \boldsymbol{\Sigma}_{z(\mathbf{r})} \mathbf{A}^t + \boldsymbol{\Sigma}_{\epsilon})$$

and

- $p(\mathbf{A}|\underline{\boldsymbol{\theta}}, \underline{\mathbf{s}}, \mathbf{z}, \underline{\mathbf{x}})$:

$$p(\mathbf{A}|\underline{\boldsymbol{\theta}}, \underline{\mathbf{s}}, \mathbf{z}, \underline{\mathbf{x}}) \propto p(\underline{\mathbf{x}}|\mathbf{A}, \underline{\mathbf{s}}, \mathbf{z}, \underline{\boldsymbol{\theta}}) p(\mathbf{A}) \quad \text{with } p(\mathbf{A}) = \mathcal{N}(0, \Gamma_p)$$

\mathbf{A} is estimated by its maximum *a posteriori*

$$\widehat{\mathbf{A}} = \arg \max_{\mathbf{A}} p(\mathbf{A}|\underline{\boldsymbol{\theta}}, \underline{\mathbf{s}}, \mathbf{z}, \underline{\mathbf{x}})$$

$$\widehat{\mathbf{A}} = \left[\sum_{\mathbf{r}} \mathbf{x}(\mathbf{r}) \bar{\mathbf{s}}'(\mathbf{r}) \right] \left[\sum_{\mathbf{r}} \bar{\mathbf{s}}(\mathbf{r}) \bar{\mathbf{s}}'(\mathbf{r}) + \mathbf{B}(\mathbf{r}) + \Gamma_p \right]^{-1}$$

5 Simulation results

For the evaluation of the results, we consider two cases: In simulation data [3], sources $\mathbf{s}(\mathbf{r})$, the common segmentation $z(\mathbf{r})$ and the mixing matrix \mathbf{A} are available. So, we can compare the estimated sources $\hat{\mathbf{s}}(\mathbf{r})$ with $\mathbf{s}(\mathbf{r})$, the estimated mixing matrix $\hat{\mathbf{A}}$ with \mathbf{A} and the estimated $\hat{z}(\mathbf{r})$ with $z(\mathbf{r})$. In the case of real data we compare $\hat{\mathbf{x}}(\mathbf{r}) = \hat{\mathbf{A}}\hat{\mathbf{s}}(\mathbf{r})$ with $\mathbf{x}(\mathbf{r})$, For comparing real valued signals $\mathbf{s}(\mathbf{r})$ and $\mathbf{x}(\mathbf{r})$, we propose to use the following Lp distances:

$\Delta_p^s = |\mathbf{s} - \hat{\mathbf{s}}|^p = \sum_{\mathbf{r}} |\mathbf{s}(\mathbf{r}) - \hat{\mathbf{s}}(\mathbf{r})|^p$ and $\Delta_p^x = |\mathbf{x} - \hat{\mathbf{x}}|^p = \sum_{\mathbf{r}} |\mathbf{x}(\mathbf{r}) - \hat{\mathbf{x}}(\mathbf{r})|^p$ for $p = 1$ and $p = 2$.

For comparing the segmentation results $z(\mathbf{r})$, we can use the number of miss-classified pixels: $\Delta^z = \sum_{\mathbf{r}} \delta(z(\mathbf{r}) - \hat{z}(\mathbf{r}))$.

When the matrix \mathbf{A} is available, then we can also use: $\Delta_p^A = |\mathbf{A} - \hat{\mathbf{A}}|^p = \sum_i \sum_j |A_{ij} - \hat{A}_{ij}|^p$. However, in most ICA methods, a separating matrix \mathbf{B} is estimated and not \mathbf{A} . Then, we use the generalized inverses $\mathbf{B}^t(\mathbf{B}\mathbf{B}^t)^{-1}$ to be compared with \mathbf{A} .

BSS Method	$\Delta_1^A(sd)$	$\Delta_2^A(sd)$	$\Delta_1^x(sd)$	$\Delta_2^x(sd)$	$\Delta_1^x(rd)$	$\Delta_2^x(rd)$
PCA	0.42	0.24	4, 18	30, 84	1, 14.10 ³	3.05.10 ³
ICA	0.45	0.27	0, 36	0, 33	1, 37.10 ³	6, 3.10 ⁶
proposed method	0.23	0.09	0, 04	0, 004	0, 11	0, 03

Table 1: Mixture matrix performances, distances between data and reconstructed data on simulated images (sd) and distances between data and reconstructed data on real images (rd)

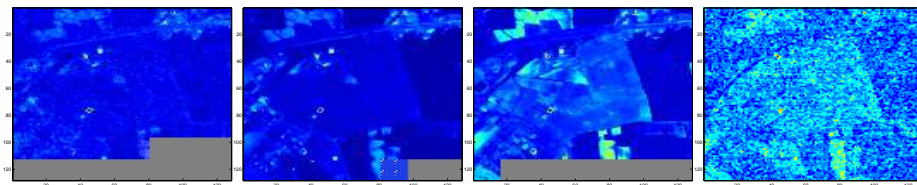


Fig. 2: 4 images extracted from 56 real AVIRIS hyperspectral data: band 1,5,30,40

6 Conclusion

In this paper, we considered the data reduction problem in hyperspectral images as a BSS and presented a Bayesian estimation approach with a particular hierarchical prior model for the observations and sources which gives us the possibility to jointly do data reduction, classification of spectra and segmentation of the images.

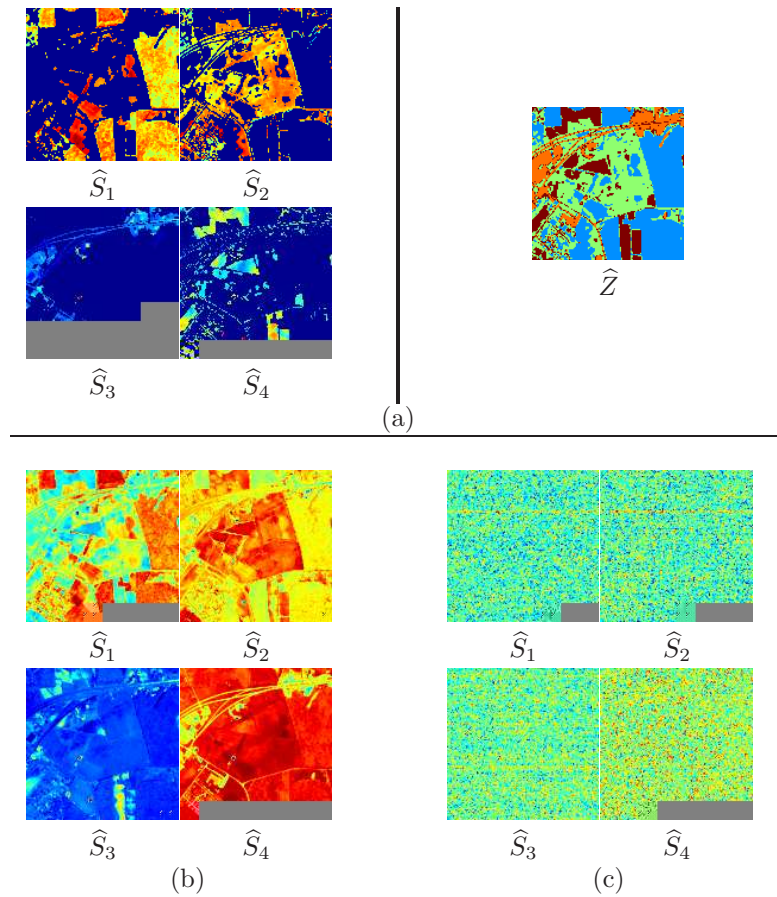


Fig. 3: results obtained size of data $(128 \times 128)\text{pixel} \times 56$ spectral band (a) proposed method reduction with segmentation (b) ICA (c) PCA. The number of sources is fixed to 4

References

- [1] K. Sasaki, S. kawata, and S. Minami, "Component analysis of spatial and spectral patterns in multispectral images," *I. Basics Journal of the Optical Society of America*, vol. 4(11), pp. 2101–2106, 1987.
- [2] C. Parra, A. Spence, A. Zieche, K.-R. Mueller, and P. Sajda, "Unmixing hyperspectral data," *In Advances in Neural Information Processing Systems*, vol. 13, pp. 848–854, 2000.
- [3] Nadia BALI and Ali Mohammad-Djafari, "Mean field approximation for bss of images with compound hierarchical gauss-markov-potts model," in *MaxEnt05*, august 2005.