



**HAL**  
open science

## CS Decomposition Based Bayesian Subspace Estimation

Olivier Besson, Nicolas Dobigeon, Jean-Yves Tournet

► **To cite this version:**

Olivier Besson, Nicolas Dobigeon, Jean-Yves Tournet. CS Decomposition Based Bayesian Subspace Estimation. IEEE Transactions on Signal Processing, 2012, 60 (8), pp.4210-4218. 10.1109/TSP.2012.2197619 . hal-03536982

**HAL Id: hal-03536982**

**<https://hal.science/hal-03536982v1>**

Submitted on 20 Jan 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



## Open Archive TOULOUSE Archive Ouverte (OATAO)

OATAO is an open access repository that collects the work of Toulouse researchers and makes it freely available over the web where possible.

This is an author-deposited version published in : <http://oatao.univ-toulouse.fr/>  
Eprints ID : 6211

To link to this article : DOI: 10.1109/ TSP.2012.2197619  
URL : <http://dx.doi.org/10.1109/TSP.2012.2197619>

To cite this version : Besson, Olivier and Dobigeon, Nicolas and Tourneret, Jean-Yves <i>CS Decomposition Based Bayesian Subspace Estimation</i> . (2012) IEEE Transactions on Signal Processing, vol. 60 (n° 8). pp. 4210-4218. ISSN 1053-587X
---

Any correspondence concerning this service should be sent to the repository administrator: [staff-oatao@listes.diff.inp-toulouse.fr](mailto:staff-oatao@listes.diff.inp-toulouse.fr)

# CS Decomposition Based Bayesian Subspace Estimation

Olivier Besson, *Senior Member, IEEE*, Nicolas Dobigeon, *Member, IEEE*, and Jean-Yves Tournet, *Senior Member, IEEE*

**Abstract**—In numerous applications, it is required to estimate the principal subspace of the data, possibly from a very limited number of samples. Additionally, it often occurs that some rough knowledge about this subspace is available and could be used to improve subspace estimation accuracy in this case. This is the problem we address herein and, in order to solve it, a Bayesian approach is proposed. The main idea consists of using the CS decomposition of the semi-orthogonal matrix  $\mathbf{H}$  whose columns span the subspace of interest. This parametrization is intuitively appealing and allows for non-informative prior distributions of the matrices involved in the CS decomposition and very mild assumptions about the angles between the actual subspace and the prior subspace. The posterior distributions are derived and a Gibbs sampling scheme is presented to obtain the minimum mean-square distance estimator of the subspace of interest. Numerical simulations and an application to real hyperspectral data assess the validity and the performances of the estimator.

**Index Terms**—Bayesian inference, CS decomposition, minimum mean-square distance estimation, simulation method, Stiefel manifold, subspace estimation.

## I. PROBLEM STATEMENT

THE ubiquitous linear model [1], [2], where the  $N$ -dimensional received signal can be written as a linear combination of  $p$  basis functions embedded in noise, has received a huge amount of attention due to its simplicity and relevance in a large number of applications. These applications include hyperspectral imagery which will be further investigated later in this paper. Under this framework, the  $N \times K$  observation matrix  $\mathbf{X}$ , where  $N$  is the dimension of the observation space and  $K$  denotes the number of measurements, can be decomposed as

$$\mathbf{X} = \mathbf{H}\Psi + \mathbf{N} \quad (1)$$

where  $\mathbf{H}$  is an  $N \times p$  matrix whose columns span the  $p$ -dimensional subspace of interest,  $\Psi$  is a  $p \times K$  matrix whose columns correspond to the coordinates of the signal in the range space  $\mathcal{R}(\mathbf{H})$  of  $\mathbf{H}$ , and  $\mathbf{N}$  denotes the additive noise. In this paper, contrary to plenty of source separation techniques such as non-negative matrix factorization or independent component analysis, we are not interested in factorizing  $\mathbf{X}$  into a product of unknown matrices  $\mathbf{H}\Psi$ . Conversely, the problem addressed in this work consists of estimating the  $p$ -dimensional subspace of interest  $\mathcal{R}(\mathbf{H})$ , which is spanned by the columns of  $\mathbf{H}$ . As a consequence, without loss of generality, we assume in the sequel that the columns of  $\mathbf{H}$  are orthonormal, i.e.,  $\mathbf{H}^T\mathbf{H} = \mathbf{I}$ . When the columns of  $\mathbf{N}$  are independent and Gaussian distributed with zero mean and covariance matrix  $\sigma^2\mathbf{I}_N$ , the maximum likelihood (ML) estimate of  $\mathcal{R}(\mathbf{H})$  is obtained from the  $p$  most significant left singular vectors of  $\mathbf{X}$  [1]. Therefore, the singular value decomposition (SVD) plays a central role in subspace estimation (in the frequentist framework) as it naturally reveals the low-rank structure of the signal. The SVD turns out to provide very accurate estimates of  $\mathcal{R}(\mathbf{H})$  in most cases [3]–[5]. However, two situations of practical interest may undermine it. The first situation corresponds to the low sample regime, a case of most interest to us as will be evidenced in the hyperspectral application of Section IV. When  $K$  is small the SVD may not produce reliable estimates: this phenomenon is especially pronounced in large dimensional problems where  $K$  might be much lower than  $N$ . In this case, the sample covariance matrix is rank-deficient and its principal subspace is poorly estimated. In order to restore a better conditioned and more accurate covariance matrix estimate, numerous techniques have been proposed including shrinkage [6], dimensionality reduction using random unitary matrices [7], constrained maximum likelihood estimation (see, e.g., [8] where the matrix of eigenvectors is constrained to be a product of Givens rotations) or eigenspace estimation using random matrix theory [9]. In the present paper, we even consider the situation where the number of snapshots  $K$  is less than the subspace dimension  $p$ . In this case, the SVD by itself is not sufficient as  $\mathbf{X}$  is at most of rank  $K < p$ , and therefore it becomes impossible to recover  $\mathcal{R}(\mathbf{H})$  without any further information. Another problem arises when the signal-to-noise ratio (SNR) is low, and hence the separation between signal singular values and noise singular values is not clear. This may result in leakage of the signal subspace into the noise subspace, or even to a subspace swap, which leads to very inaccurate subspace estimates. This phenomenon has been evidenced, e.g., in [10] and [11], and theoretical explanations, based on the theory of large dimensional random matrices [12]

are now available to predict this behavior [13]–[15]. In fact, for the two cases mentioned previously, additional prior information may prove to be helpful, and this prior information is often available either through models, expertise or data (cf. the hyperspectral application studied later in this paper). A natural way to introduce such knowledge is to adhere to a Bayesian framework. This is the approach we advocate in the present paper where our main focus is on *knowledge-aided subspace estimation in the low sample support or low SNR regime*.

More precisely, we assume that  $\mathbf{H}$  is assigned some prior distribution  $\pi(\mathbf{H})$ , and our goal is to estimate  $\mathbf{H}$  from the posterior distribution  $p(\mathbf{H}|\mathbf{X})$ . Similarly to [16] and [17], we consider minimum mean-square distance (MMSD) estimators of  $\mathbf{H}$ , i.e., we look for estimates  $\hat{\mathbf{H}}$  of  $\mathbf{H}$  that minimize the average squared Frobenius norm of the difference between the projection matrices, viz.,  $E\{\|\hat{\mathbf{H}}\hat{\mathbf{H}}^T - \mathbf{H}\mathbf{H}^T\|_F^2\}$ . The rationale behind this approach is that the usual mean-square metric  $E\{\|\hat{\mathbf{H}} - \mathbf{H}\|_F^2\}$  is not the natural metric on the Stiefel manifold [18], [19] while the distance between projection matrices is meaningful.<sup>1</sup> Using the latter distance, the MMSD estimator was shown to be given by [16], [17]

$$\begin{aligned}\hat{\mathbf{H}}_{\text{mmsd}} &= \arg \max_{\hat{\mathbf{H}}} E \left\{ \text{Tr} \left\{ \hat{\mathbf{H}}^T \mathbf{H} \mathbf{H}^T \hat{\mathbf{H}} \right\} \right\} \\ &= \arg \max_{\hat{\mathbf{H}}} \int \text{Tr} \left\{ \hat{\mathbf{H}}^T \mathbf{H} \mathbf{H}^T \hat{\mathbf{H}} \right\} p(\mathbf{H}|\mathbf{X}) d\mathbf{H} \\ &= \arg \max_{\hat{\mathbf{H}}} \text{Tr} \left\{ \hat{\mathbf{H}}^T \left[ \int \mathbf{H} \mathbf{H}^T p(\mathbf{H}|\mathbf{X}) d\mathbf{H} \right] \hat{\mathbf{H}} \right\} \\ &= \mathcal{P}_p \left\{ \int \mathbf{H} \mathbf{H}^T p(\mathbf{H}|\mathbf{X}) d\mathbf{H} \right\}\end{aligned}\quad (2)$$

where  $\mathcal{P}_p\{\cdot\}$  stands for the  $p$  principal eigenvectors of the matrix between braces. The MMSD estimator thus amounts to finding the principal subspace of the posterior mean of the projection matrix  $\mathbf{P} = \mathbf{H}\mathbf{H}^T$  on  $\mathcal{R}(\mathbf{H})$ . Note that this approach is general and independent of the conditional and prior distributions: depending on the latter, it may or may not be an easy task to obtain the MMSD estimator. In the sequel, we state our assumptions regarding  $\mathbf{H}$  and derive its corresponding MMSD estimator. The latter will then be tested on real hyperspectral data in Section IV.

## II. DATA MODEL AND SUBSPACE ESTIMATION

Let us consider the linear model (1) and let us assume that  $\mathbf{N}$  is Gaussian distributed with independent columns so that the probability density function of  $\mathbf{X}$ , conditioned on  $\mathbf{H}$  and  $\Psi$ , is given by

$$p(\mathbf{X}|\mathbf{H}, \Psi) \propto \text{etr} \left\{ -\frac{1}{2\sigma^2} (\mathbf{X} - \mathbf{H}\Psi)^T (\mathbf{X} - \mathbf{H}\Psi) \right\} \quad (3)$$

<sup>1</sup>The true (square) distance between the subspaces is given by  $d^2(\hat{\mathbf{H}}, \mathbf{H}) = \sum_{k=1}^p \theta_k^2$ , where  $\theta_k$  for  $k = 1, \dots, p$  stand for the principal angles between  $\mathcal{R}(\hat{\mathbf{H}})$  and  $\mathcal{R}(\mathbf{H})$ . The distance we use herein, i.e.,  $\|\hat{\mathbf{H}}\hat{\mathbf{H}}^T - \mathbf{H}\mathbf{H}^T\|_F^2 = 2p - 2\text{Tr}\{\hat{\mathbf{H}}^T \mathbf{H} \mathbf{H}^T \hat{\mathbf{H}}\} = 2 \sum_{k=1}^p \sin^2 \theta_k$ , is thus different from  $d^2(\hat{\mathbf{H}}, \mathbf{H})$ . However, the two distances are close for small values of  $\theta_k$  and the distance between projection matrices is widely accepted. Moreover, using the distance between projection matrices allows one to obtain a closed-form expression for the MMSD estimator, see (2). Minimization of  $E\{d^2(\hat{\mathbf{H}}, \mathbf{H})\}$  would not yield such closed-form expression since  $\sum_{k=1}^p \theta_k^2$  cannot be expressed simply as a function of  $\hat{\mathbf{H}}$  and  $\mathbf{H}$ .

where  $\text{etr}\{\cdot\}$  stands for the exponential of the trace of the matrix between braces and  $\propto$  means proportional to. Since the thermal noise level can usually be estimated with high accuracy, we assume here that  $\sigma^2$  is known.<sup>2</sup> Since no knowledge about  $\Psi$  is generally available, we treat it as a random matrix with uniform prior distribution, i.e.,  $\pi(\Psi) \propto 1$ , so that the distribution of  $\mathbf{X}$ , conditioned on  $\mathbf{H}$  only, is obtained as

$$\begin{aligned}p(\mathbf{X}|\mathbf{H}) &= \int p(\mathbf{X}|\mathbf{H}, \Psi) \pi(\Psi) d\Psi \\ &\propto \int \text{etr} \left\{ -\frac{1}{2\sigma^2} (\mathbf{X} - \mathbf{H}\Psi)^T (\mathbf{X} - \mathbf{H}\Psi) \right\} d\Psi \\ &\propto \text{etr} \left\{ -\frac{1}{2\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{2\sigma^2} \mathbf{X}^T \mathbf{H} \mathbf{H}^T \mathbf{X} \right\} \\ &\quad \times \int \text{etr} \left\{ -\frac{1}{2\sigma^2} (\Psi - \mathbf{H}^T \mathbf{X})^T (\Psi - \mathbf{H}^T \mathbf{X}) \right\} d\Psi \\ &\propto \text{etr} \left\{ -\frac{1}{2\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{2\sigma^2} \mathbf{X}^T \mathbf{H} \mathbf{H}^T \mathbf{X} \right\}\end{aligned}\quad (4)$$

where, to obtain the last line, we have used the fact that the integral in the fourth line of (4) is that of a multivariate Gaussian distribution with mean  $\mathbf{H}^T \mathbf{X}$  and covariance matrix  $\sigma^2 \mathbf{I}$ , and hence is proportional to  $\sigma^{pK}$ . Note that  $p(\mathbf{X}|\mathbf{H})$  depends on  $\mathbf{H}$  only through the projection matrix  $\mathbf{P} = \mathbf{H}\mathbf{H}^T$ .

Let us turn now to the hypotheses regarding  $\mathbf{H}$ . We assume that we have some *a priori* knowledge about the subspace spanned by the columns of  $\mathbf{H}$ : This knowledge can come from some available models or can be deduced from the data itself, as in the hyperspectral imagery application. More precisely, we assume that the range space  $\mathcal{R}(\mathbf{H})$  of  $\mathbf{H}$  is close to the range space of some semi-orthogonal matrix  $\bar{\mathbf{H}}$  and, without loss of generality, we will assume that  $\bar{\mathbf{H}} = [\mathbf{I}_p \ \mathbf{0}]^T$  through the paper.<sup>3</sup>

In [17], we tackled the problem by assigning the matrix  $\mathbf{H}$  either a Bingham— $\pi_B(\mathbf{H}) \propto \text{etr}\{\kappa \mathbf{H}^T \bar{\mathbf{H}} \bar{\mathbf{H}}^T \mathbf{H}\}$ —or a von Mises–Fisher (vMF) distribution— $\pi_{\text{vMF}}(\mathbf{H}) \propto \text{etr}\{\kappa \mathbf{H}^T \bar{\mathbf{H}}\}$ . The Bingham and vMF are the most widely used distributions on the Stiefel manifold and they have proved to be relevant in a number of applications, including meteorology, biology, image, or shape analysis [20]. Moreover, there exists computationally efficient simulation tools to sample from these distributions, which makes them a sensible choice. However, they suffer from two drawbacks. First, from a user point of view, it is not obvious to set a value for the concentration parameter  $\kappa$  since the latter is not an intuitively appealing parameter, in contrast to the angles between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$  which are more directly meaningful. Moreover, the Bingham and vMF distributions hold for the whole matrix  $\mathbf{H}$ : the choice

<sup>2</sup>The case of unknown  $\sigma^2$  can be considered by assigning a prior distribution (typically a conjugate prior, in our case an inverse gamma distribution) to  $\sigma^2$  and modifying accordingly the posterior distributions to be derived next.

<sup>3</sup>In the case where  $\mathbf{H}$  is close to an arbitrary semi-orthogonal matrix  $\bar{\mathbf{H}}$ , the measurements in (1) can be pre-multiplied by the unitary matrix  $\mathbf{Q}$  such that  $\mathbf{Q}\bar{\mathbf{H}} = [\mathbf{I}_p \ \mathbf{0}]^T$ . Note that pre-multiplication by the unitary matrix  $\mathbf{Q}$  does not modify the angles between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$  nor the distribution  $p(\mathbf{X}|\mathbf{H}, \Psi)$  in (3).

of a distribution and a value for  $\kappa$  will consequently induce a distribution for the angles, but this relation is not revealed in a straightforward and intelligible manner. In the present paper, we attempt to remedy these shortcomings with a view to obtain a parametrization of the statistical model that directly involves the most meaningful parameters, namely the angles  $\theta_k$ ,  $k = 1, \dots, p$  between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$ . Indeed, these angles are instrumental as the distance between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$  is directly connected to them. Furthermore, we look for a less constrained model which relies on mild assumptions, and the latter would only concern the angles  $\theta_k$ .

The model proposed herein is based on the CS decomposition of  $\mathbf{H}$ , which writes [19]

$$\mathbf{H} = \begin{bmatrix} \mathbf{U}_1 \mathbf{C} \\ \mathbf{U}_2 \mathbf{S} \end{bmatrix} \mathbf{V}^T \quad (5)$$

where  $\mathbf{U}_1$  and  $\mathbf{V}$  are  $p \times p$  orthogonal matrices,  $\mathbf{U}_2$  is an  $(N - p) \times p$  semi-orthogonal matrix ( $\mathbf{U}_2^T \mathbf{U}_2 = \mathbf{I}_p$ ),  $\mathbf{C} = \text{diag}(\cos \theta_1, \dots, \cos \theta_p)$  and  $\mathbf{S} = \text{diag}(\sin \theta_1, \dots, \sin \theta_p)$ . The angles  $\theta_k$  correspond to the principal angles between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$  while the columns of  $\begin{bmatrix} \mathbf{U}_1 \\ \mathbf{0} \end{bmatrix}$  and  $\mathbf{H}\mathbf{V}$  are the associated principal vectors. As requested, this representation has the nice property that the angles between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$  are directly revealed, and do not depend on the matrices  $\mathbf{U}_1$ ,  $\mathbf{U}_2$ , and  $\mathbf{V}$ , which can be arbitrary. We now assign prior distributions to the model variables. First observe that the likelihood function in (4) depends on  $\mathbf{H}$  only through the projection matrix  $\mathbf{P} = \mathbf{H}\mathbf{H}^T$  and the latter, under the CS decomposition (5), is independent of  $\mathbf{V}$ . Therefore, we need to set prior distributions for  $\mathbf{U}_1$ ,  $\mathbf{U}_2$ , and  $\boldsymbol{\theta} = [\theta_1 \ \dots \ \theta_p]^T$  only. As for  $\mathbf{U}_1$  and  $\mathbf{U}_2$ , we assume that they have uniform prior distributions on the orthogonal group  $\mathcal{O}(p)$  and the Stiefel manifold  $\mathcal{S}_{p, N-p}$ , i.e., the set of  $(N - p) \times p$  matrices  $\mathbf{U}_2$  such that  $\mathbf{U}_2^T \mathbf{U}_2 = \mathbf{I}_p$ . As for  $\boldsymbol{\theta}$ , we assume that  $\theta_k$  are independent and identically distributed (i.i.d.) random variables, with uniform distribution on  $[0, \theta_{\max}]$ , i.e.,  $\theta_k \sim U([0, \theta_{\max}])$ . Observe that, as stated in our objectives, the statistical model involves rather mild assumptions. Moreover, it directly involves the angles  $\theta_k$ , which makes sense intuitively. Finally, the only parameter the user has to set is  $\theta_{\max}$ , which seems easier to set than a value for  $\kappa$ . Indeed  $\theta_{\max}$  rules the maximum angle between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$ : therefore, the smaller  $\theta_{\max}$ , the closer these subspaces *a priori*. In contrast, when  $\theta_{\max}$  increases, the two subspaces can be quite far apart. Consequently, for small  $\theta_{\max}$ , we can expect the MMSD estimator to strongly rely on  $\bar{\mathbf{H}}$ , while for large  $\theta_{\max}$  the data  $\mathbf{X}$  is likely to prevail.

Since the likelihood and the prior distributions have been set, we now consider the posterior distributions of  $\mathbf{U}_1$ ,  $\mathbf{U}_2$ , and  $\boldsymbol{\theta}$ . As a preliminary step, note that

$$\mathbf{P} = \mathbf{H}\mathbf{H}^T = \begin{bmatrix} \mathbf{U}_1 \mathbf{C}^2 \mathbf{U}_1^T & \mathbf{U}_1 \mathbf{C} \mathbf{S} \mathbf{U}_2^T \\ \mathbf{U}_2 \mathbf{S} \mathbf{C} \mathbf{U}_1^T & \mathbf{U}_2 \mathbf{S}^2 \mathbf{U}_2^T \end{bmatrix} \quad (6)$$

so that, with the partitioning  $\mathbf{X} = [\mathbf{X}_1^T \ \mathbf{X}_2^T]^T$ , we have

$$\begin{aligned} & \text{Tr} \{ \mathbf{X}^T \mathbf{P} \mathbf{X} \} \\ &= \text{Tr} \left\{ \begin{bmatrix} \mathbf{X}_1^T & \mathbf{X}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{U}_1 \mathbf{C}^2 \mathbf{U}_1^T & \mathbf{U}_1 \mathbf{C} \mathbf{S} \mathbf{U}_2^T \\ \mathbf{U}_2 \mathbf{S} \mathbf{C} \mathbf{U}_1^T & \mathbf{U}_2 \mathbf{S}^2 \mathbf{U}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} \right\} \\ &= \text{Tr} \left\{ \mathbf{C}^2 \mathbf{U}_1^T \mathbf{X}_1 \mathbf{X}_1^T \mathbf{U}_1 + 2 \mathbf{X}_2^T \mathbf{U}_2 \mathbf{S} \mathbf{C} \mathbf{U}_1^T \mathbf{X}_1 \right. \\ & \quad \left. + \mathbf{S}^2 \mathbf{U}_2^T \mathbf{X}_2 \mathbf{X}_2^T \mathbf{U}_2 \right\}. \end{aligned} \quad (7)$$

Assuming *a priori* independence between  $\mathbf{U}_1$  and  $\mathbf{U}_2$  and  $\boldsymbol{\theta}$ , it follows from (4) that the joint posterior distribution of  $\mathbf{U}_1$  and  $\mathbf{U}_2$  and  $\boldsymbol{\theta}$  is given by

$$\begin{aligned} & p(\mathbf{U}_1, \mathbf{U}_2, \boldsymbol{\theta} | \mathbf{X}) \\ & \propto p(\mathbf{X} | \mathbf{H}) \pi(\mathbf{U}_1) \pi(\mathbf{U}_2) \pi(\boldsymbol{\theta}) \\ & \propto \text{etr} \left\{ \frac{1}{2\sigma^2} [\mathbf{C}^2 \mathbf{U}_1^T \mathbf{X}_1 \mathbf{X}_1^T \mathbf{U}_1 + \mathbf{S}^2 \mathbf{U}_2^T \mathbf{X}_2 \mathbf{X}_2^T \mathbf{U}_2] \right\} \\ & \quad \times \text{etr} \left\{ \frac{1}{\sigma^2} \mathbf{X}_2^T \mathbf{U}_2 \mathbf{S} \mathbf{C} \mathbf{U}_1^T \mathbf{X}_1 \right\} \pi(\mathbf{U}_1) \pi(\mathbf{U}_2) \pi(\boldsymbol{\theta}). \end{aligned} \quad (8)$$

In order to obtain the MMSD estimator, we suggest, as in [17], to use a Gibbs sampler which enables one to iteratively draw samples from the posterior distribution of each variable, conditioned on all other variables [21], [22]. In order to obtain the conditional posterior distribution of  $\mathbf{U}_1$  only, we start with (8) and keep only the terms which depend on  $\mathbf{U}_1$  since the other terms will appear as constants and can be absorbed in the normalization constant. Doing so, we deduce that

$$\begin{aligned} & p(\mathbf{U}_1 | \mathbf{U}_2, \boldsymbol{\theta}, \mathbf{X}) \propto \text{etr} \left\{ \frac{1}{2\sigma^2} \mathbf{C}^2 \mathbf{U}_1^T \mathbf{X}_1 \mathbf{X}_1^T \mathbf{U}_1 \right\} \\ & \quad \times \text{etr} \left\{ \frac{1}{\sigma^2} \mathbf{U}_1^T \mathbf{X}_1 \mathbf{X}_2^T \mathbf{U}_2 \mathbf{S} \mathbf{C} \right\} \mathbb{1}_{\mathcal{O}(p)}(\mathbf{U}_1) \end{aligned} \quad (9)$$

where  $\mathbb{1}_{\mathcal{O}(p)}(\mathbf{U}_1)$  is the indicator function defined on  $\mathcal{O}(p)$  (i.e.,  $\mathbb{1}_{\mathcal{O}(p)}(\mathbf{U}_1) = 1$  if  $\mathbf{U}_1 \in \mathcal{O}(p)$  and 0 otherwise). The distribution in (9) is recognized as a Bingham–von Mises–Fisher (BMF) distribution with parameter matrices  $\mathbf{X}_1 \mathbf{X}_1^T$ ,  $\frac{1}{2\sigma^2} \mathbf{C}^2$  and  $\frac{1}{\sigma^2} \mathbf{X}_1 \mathbf{X}_2^T \mathbf{U}_2 \mathbf{S} \mathbf{C}$  respectively.<sup>4</sup> An efficient sampling scheme to generate random matrices drawn from a BMF  $(\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3)$  distribution on the Stiefel manifold was proposed in [23]. In our case,  $\mathbf{U}_1 \in \mathcal{O}(p)$  and, as mentioned in [23], the sampling scheme on the Stiefel manifold cannot be used directly and needs to be modified. In [24, App. A], following the lines of [23], we give some details about the sampling scheme for a matrix BMF distribution on the orthogonal group  $\mathcal{O}(p)$ . Similarly, we have

$$\begin{aligned} & p(\mathbf{U}_2 | \mathbf{U}_1, \boldsymbol{\theta}, \mathbf{X}) \propto \text{etr} \left\{ \frac{1}{2\sigma^2} \mathbf{S}^2 \mathbf{U}_2^T \mathbf{X}_2 \mathbf{X}_2^T \mathbf{U}_2 \right\} \\ & \quad \times \text{etr} \left\{ \frac{1}{\sigma^2} \mathbf{U}_2^T \mathbf{X}_2 \mathbf{X}_1^T \mathbf{U}_1 \mathbf{C} \mathbf{S} \right\} \mathbb{1}_{\mathcal{S}_{p, N-p}}(\mathbf{U}_2) \end{aligned} \quad (10)$$

<sup>4</sup>The matrix  $\mathbf{H} \in \mathcal{S}_{p, q}$  is said to have a BMF  $(\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3)$  distribution where  $\mathbf{A}_1$  is an  $q \times q$  symmetric matrix,  $\mathbf{A}_2$  is a  $p \times p$  diagonal matrix and  $\mathbf{A}_3$  is an  $q \times p$  matrix- if  $p(\mathbf{H}) \propto \text{etr} \{ \mathbf{H}^T \mathbf{A}_3 + \mathbf{A}_2 \mathbf{H}^T \mathbf{A}_1 \mathbf{H} \}$ .

TABLE I  
GIBBS SAMPLER FOR ESTIMATION OF  $\mathbf{H}$  USING THE CS DECOMPOSITION

**Input:** initial values  $\mathbf{U}_1^{(0)}, \mathbf{U}_2^{(0)}, \boldsymbol{\theta}^{(0)}$   
1: **for**  $n = 1, \dots, N_{\text{bi}} + N_r$  **do**  
2: sample  $\mathbf{U}_1^{(n)}$  from BMF  $\left(\mathbf{X}_1 \mathbf{X}_1^T, \frac{1}{2\sigma^2} (\mathbf{C}^{(n-1)})^2, \frac{1}{\sigma^2} \mathbf{X}_1 \mathbf{X}_2^T \mathbf{U}_2^{(n-1)} \mathbf{S}^{(n-1)} \mathbf{C}^{(n-1)}\right)$  in (9) where  $\mathbf{C}^{(n-1)} = \text{diag}(\cos \boldsymbol{\theta}^{(n-1)})$  and  $\mathbf{S}^{(n-1)} = \text{diag}(\sin \boldsymbol{\theta}^{(n-1)})$ .  
3: sample  $\mathbf{U}_2^{(n)}$  from BMF  $\left(\mathbf{X}_2 \mathbf{X}_2^T, \frac{1}{2\sigma^2} (\mathbf{S}^{(n-1)})^2, \frac{1}{\sigma^2} \mathbf{X}_2 \mathbf{X}_1^T \mathbf{U}_1^{(n)} \mathbf{C}^{(n-1)} \mathbf{S}^{(n-1)}\right)$  in (10).  
4: **Metropolis-Hastings** to sample  $\boldsymbol{\theta}^{(n)}$ :  
5: **for**  $k = 1, \dots, p$  **do**  
6: draw an initial candidate  $x_k^{c(0)}$  from  $x_{\max} \times \text{Beta}(a_k, b_k)$  and set  $x_k^{(n)} = x_k^{c(0)}$ .  
7: **for**  $\ell = 1, \dots, q$  **do**  
8: draw a candidate  $x_k^{c(\ell)}$  from  $x_{\max} \times \text{Beta}(a_k, b_k)$ .  
9: accept  $x_k^{c(\ell)}$  as  $x_k^{(n)}$  with probability  $\min\left(1, \frac{p(x_k^{c(\ell)} | \mathbf{U}_1^{(n)}, \mathbf{U}_2^{(n)}, \mathbf{X})}{p(x_k^{c(\ell-1)} | \mathbf{U}_1^{(n)}, \mathbf{U}_2^{(n)}, \mathbf{X})} \frac{q(x_k^{c(\ell-1)})}{q(x_k^{c(\ell)})}\right)$  where  $p(x_k | \mathbf{U}_1^{(n)}, \mathbf{U}_2^{(n)}, \mathbf{X})$  is given in (13).  
10: **end for**  
11:  $\boldsymbol{\theta}_k^{(n)} = \arcsin \sqrt{x_k^{(n)}}$ .  
12: **end for**  
13: **end for**  
**Output:** sequence of random matrices  $\mathbf{H}^{(n)} = \begin{bmatrix} \mathbf{U}_1^{(n)} \mathbf{C}^{(n)} \\ \mathbf{U}_2^{(n)} \mathbf{S}^{(n)} \end{bmatrix}$ .

and hence

$$\mathbf{U}_2 | \mathbf{U}_1, \boldsymbol{\theta}, \mathbf{X} \sim \text{BMF} \left( \mathbf{X}_2 \mathbf{X}_2^T, \frac{1}{2\sigma^2} \mathbf{S}^2, \frac{1}{\sigma^2} \mathbf{X}_2 \mathbf{X}_1^T \mathbf{U}_1 \mathbf{C} \mathbf{S} \right). \quad (11)$$

Since  $\mathbf{U}_2 \in \mathcal{S}_{p, N-p}$ , the sampling scheme of Hoff [23] can be used to draw matrices from the distribution in (11). Let us now examine the posterior distribution of  $\boldsymbol{\theta}$

$$\begin{aligned} p(\boldsymbol{\theta} | \mathbf{U}_1, \mathbf{U}_2, \mathbf{X}) &\propto \text{etr} \left\{ \frac{1}{\sigma^2} \mathbf{X}_2^T \mathbf{U}_2 \mathbf{S} \mathbf{C} \mathbf{U}_1^T \mathbf{X}_1 \right\} \\ &\times \text{etr} \left\{ \frac{1}{2\sigma^2} [\mathbf{C}^2 \mathbf{U}_1^T \mathbf{X}_1 \mathbf{X}_1^T \mathbf{U}_1 + \mathbf{S}^2 \mathbf{U}_2^T \mathbf{X}_2 \mathbf{X}_2^T \mathbf{U}_2] \right\} \pi(\boldsymbol{\theta}) \\ &\propto \prod_{k=1}^p e^{\alpha_k \cos^2 \theta_k + 2\beta_k \cos \theta_k \sin \theta_k + \gamma_k \sin^2 \theta_k} \mathbb{I}_{[0, \theta_{\max}]}(\theta_k) \quad (12) \end{aligned}$$

where  $\alpha_k, \beta_k, \gamma_k$  are the  $k$ th diagonal entries of  $\frac{1}{2\sigma^2} \mathbf{U}_1^T \mathbf{X}_1 \mathbf{X}_1^T \mathbf{U}_1$ ,  $\frac{1}{2\sigma^2} \mathbf{U}_1^T \mathbf{X}_1 \mathbf{X}_2^T \mathbf{U}_2$  and  $\frac{1}{2\sigma^2} \mathbf{U}_2^T \mathbf{X}_2 \mathbf{X}_2^T \mathbf{U}_2$ , respectively. The first thing to be noted is that the variables  $\theta_k$ , conditioned on  $\mathbf{U}_1, \mathbf{U}_2$  and  $\mathbf{X}$ , are independent and hence one needs to generate  $p$  independent random variables. Unfortunately, the distribution in (12) does not belong to any known class of distributions and, therefore, generating random variables drawn from  $p(\boldsymbol{\theta} | \mathbf{U}_1, \mathbf{U}_2, \mathbf{X})$  appears problematic. In order to overcome this problem, we propose to resort to a Metropolis–Hastings (MH) move [21], [22]. The basic idea is to generate a random variable drawn from a proposal distribution and to accept it with a certain probability, the latter being equal to one if the candidate contributes to increase the target posterior distribution. Of course, the closer the proposal and target distributions, the higher the acceptance rate and hence the faster the convergence of the Markov chain. In order to obtain a proposal distribution in our case, we make the change of variable  $x_k = \sin^2 \theta_k$  in (12), and come up with the

equivalent problem of finding a proposal distribution for the conditional distribution of  $x_k$ , which is given by

$$p(x_k | \mathbf{U}_1, \mathbf{U}_2, \mathbf{X}) \propto e^{-(\alpha_k - \gamma_k)x_k + 2\beta_k x_k^{\frac{1}{2}} (1-x_k)^{\frac{1}{2}}} \times x_k^{\frac{1}{2}} (1-x_k)^{\frac{1}{2}} \mathbb{I}_{[0, x_{\max}]}(x_k) \quad (13)$$

where  $x_{\max} = \sin^2 \theta_{\max}$ . Forgetting the exponential term in (13), this distribution is similar to that of a scaled beta distribution. Therefore, we choose a scaled beta distribution  $q(x_k) \propto \left(\frac{x_k}{x_{\max}}\right)^{a_k-1} \left(1 - \frac{x_k}{x_{\max}}\right)^{b_k-1}$  as a proposal distribution in a Metropolis–Hastings scheme. Through preliminary investigation, we ended up with the choice  $a_k = 0.5 + 0.25 \max(0, \beta_k) - 0.25 \min(0, \alpha_k - \gamma_k)$  and  $b_k = 0.5 + 0.25 \max(0, \beta_k) + 0.25 \max(0, \alpha_k - \gamma_k)$  which turns out to provide a good approximation to (13) for low to moderate SNR. The resulting Gibbs sampling scheme is summarized in Table I.

Once the matrices  $\mathbf{H}^{(n)}$  have been generated, the MMSD estimator, which theoretically entails computing  $\mathcal{P}_p \left\{ \int \mathbf{H} \mathbf{H}^T p(\mathbf{H} | \mathbf{X}) d\mathbf{H} \right\}$ , can be approximated by

$$\hat{\mathbf{H}}_{\text{mmsd}} = \mathcal{P}_p \left\{ \frac{1}{N_r} \sum_{n=N_{\text{bi}}+1}^{N_{\text{bi}}+N_r} \mathbf{H}^{(n)} (\mathbf{H}^{(n)})^T \right\}. \quad (14)$$

*Remark 1:* Similarly, a maximum a posteriori (MAP) approach can be advocated where the MAP estimator is obtained as

$$\begin{aligned} \hat{\mathbf{H}}_{\text{map}} &= \arg \max_{\mathbf{H}^{(n)}} p \left( \mathbf{U}_1^{(n)}, \mathbf{U}_2^{(n)}, \boldsymbol{\theta}^{(n)} | \mathbf{X} \right) \\ &= \arg \max_{\mathbf{H}^{(n)}} \text{Tr} \left\{ \mathbf{X}^T \mathbf{H}^{(n)} (\mathbf{H}^{(n)})^T \mathbf{X} \right\}. \quad (15) \end{aligned}$$

Note that  $\text{Tr} \{ \mathbf{X}^T \mathbf{H} \mathbf{H}^T \mathbf{X} \}$  is maximized when  $\mathbf{H}$  is the matrix of the  $p$  most significant left singular vectors of  $\mathbf{X}$  and, hence, the MAP approach is in some way linked to the SVD-based approach. Observe also that it does not make much sense to consider here a minimum mean-square error (MMSE) estimator. Indeed the latter entails computing  $\int \mathbf{H} p(\mathbf{H} | \mathbf{X}) d\mathbf{H}$ , which could be approximated by the arithmetic mean of the set of matrices

$\mathbf{H}^{(n)}$ . However, the range space of  $\mathbf{H}$  is given up to right multiplication by an orthogonal matrix. Therefore,  $\mathcal{R}(\mathbf{H}^{(n)})$  could be close to  $\mathcal{R}(\mathbf{H})$  without the actual matrices  $\mathbf{H}^{(n)}$  and  $\mathbf{H}$  being close. It results that the arithmetic mean of the matrices  $\mathbf{H}^{(n)}$  could result in a poor subspace estimate despite the fact that, individually, the subspaces spanned by each matrix  $\mathbf{H}^{(n)}$  might be accurate.

### III. SIMULATIONS

In this section, we use Monte Carlo simulations to assess the performance of the estimator defined previously. The performance measure will be the distance between the subspace spanned by  $\hat{\mathbf{H}}$  and the subspace spanned by  $\mathbf{H}$  where  $\hat{\mathbf{H}}$  stands for one of the estimates. More precisely, we will display the mean-square distance (MSD), which is defined as

$$\text{MSD}(\hat{\mathbf{H}}, \mathbf{H}) = \mathbb{E}\{d^2(\hat{\mathbf{H}}, \mathbf{H})\} = \mathbb{E}\left\{\sum_{k=1}^p \theta_k^2\right\} \quad (16)$$

where  $\theta_k$ ,  $k = 1, \dots, p$  stand for the principal angles between  $\mathcal{R}(\hat{\mathbf{H}})$  and  $\mathcal{R}(\mathbf{H})$ . In all simulations,  $N = 20$ ,  $p = 5$ , and  $\hat{\mathbf{H}} = [\mathbf{I}_p \ \mathbf{0}]^T$ . The matrix  $\Psi$  is generated from a Gaussian distribution with zero-mean and covariance matrix  $\mathbf{I}_p$ , and the SNR is defined as

$$\text{SNR} = 10 \log_{10} \left( \frac{\mathbb{E}\{\text{Tr}\{\Psi^T \mathbf{H}^T \mathbf{H} \Psi\}\}}{\mathbb{E}\{\text{Tr}\{\mathbf{N}^T \mathbf{N}\}\}} \right) = 10 \log_{10} \left( \frac{p}{N\sigma^2} \right).$$

The angles between  $\mathcal{R}(\mathbf{H})$  and  $\mathcal{R}(\bar{\mathbf{H}})$  are fixed over all simulations and set to  $\boldsymbol{\theta} = [15^\circ \ 25^\circ \ 35^\circ \ 45^\circ \ 55^\circ]^T$ , which results in  $\text{MSD}(\bar{\mathbf{H}}, \mathbf{H}) = 2.1704$ . The matrices  $\mathbf{U}_1$  and  $\mathbf{U}_2$  are drawn randomly at each Monte Carlo run. The number of burn-in iterations in the Gibbs sampler is set to  $N_{\text{bi}} = 10$  and  $N_r = 1000$  samples are used to approximate the estimators following (14) and (15). The MMSD estimator(14) is compared with the usual SVD-based estimator and the sparse matrix transform (SMT) of [8]. In all figures, the solid black line represents  $\text{MSD}(\bar{\mathbf{H}}, \mathbf{H})$ , i.e., when  $\hat{\mathbf{H}} = \bar{\mathbf{H}}$  and only the *a priori* knowledge is used, the data being discarded. In all simulations, the MSD is evaluated from 500 Monte Carlo trials. Additional simulation results, including a comparison with the MAP estimator and evaluation of the average fraction of energy of  $\hat{\mathbf{H}}$  in  $\mathcal{R}(\mathbf{H})$ , can be found in [24].

We now successively investigate the influence of  $\theta_{\max}$ ,  $K$  and SNR in Figs. 1–9. The first observation to be made is that the MMSD estimator is rather insensitive to the choice of  $\theta_{\max}$ : this is an interesting feature, as it means that  $\theta_{\max}$  need not be fixed with a high accuracy. Next, it can be observed that the MMSD estimator outperforms the usual SVD-based estimator and the SMT estimator, for small  $K$  and low SNR: Under these conditions, it makes a sound use of the prior information and provides more accurate estimates. Note also that it performs better than the estimate  $\hat{\mathbf{H}} = \bar{\mathbf{H}}$ , and hence the prior by itself is not sufficient. Finally, we observe that SMT is approximately equivalent to the SVD estimator.

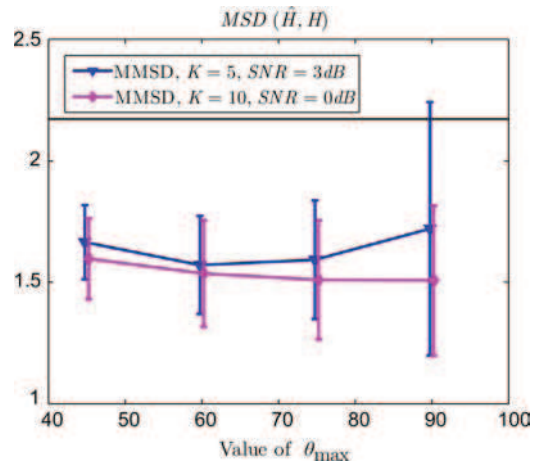


Fig. 1. Mean-square distance between true and estimated subspaces versus  $\theta_{\max}$ .  $N = 20$ ,  $p = 5$ .

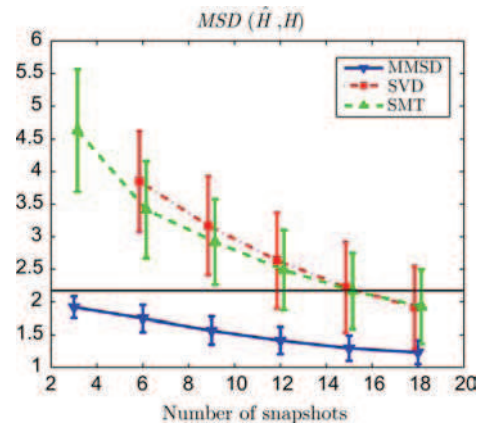


Fig. 2. Mean-square distance between true and estimated subspaces versus  $K$ .  $N = 20$ ,  $p = 5$ , SNR = 0 dB and  $\theta_{\max} = 60^\circ$ .

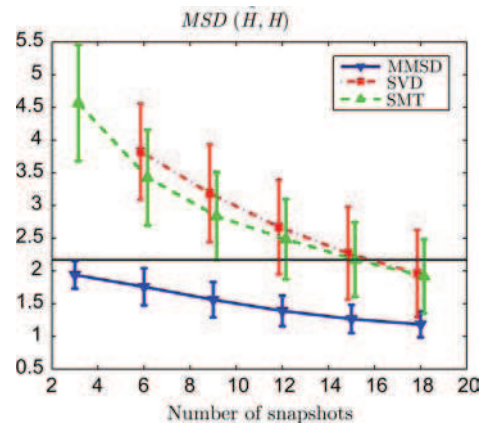


Fig. 3. Mean-square distance between true and estimated subspaces versus  $K$ .  $N = 20$ ,  $p = 5$ , SNR = 0 dB, and  $\theta_{\max} = 75^\circ$ .

### IV. APPLICATION TO HYPERSPECTRAL DATA

Hyperspectral imagery has recently emerged as a feasible and relevant technique for accurate observation of earth surfaces, either for agricultural or geographical purposes [25]. The diversity of the frequency response of each component of the illuminated scene makes it possible to gain a fine understanding

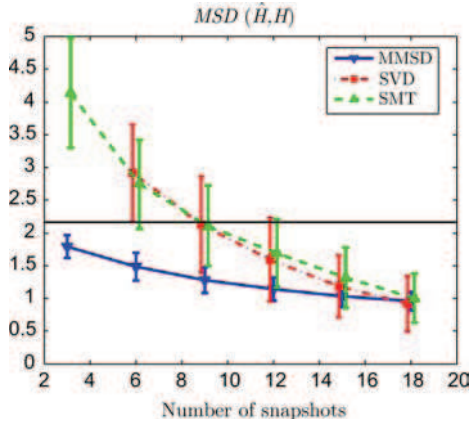


Fig. 4. Mean-square distance between true and estimated subspaces versus  $K$ .  $N = 20$ ,  $p = 5$ , SNR = 3 dB and  $\theta_{\max} = 60^\circ$ .

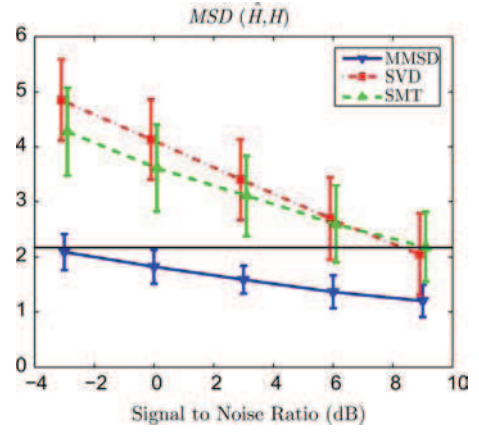


Fig. 7. Mean-square distance between true and estimated subspaces versus SNR.  $N = 20$ ,  $p = 5$ ,  $K = 5$ , and  $\theta_{\max} = 75^\circ$ .

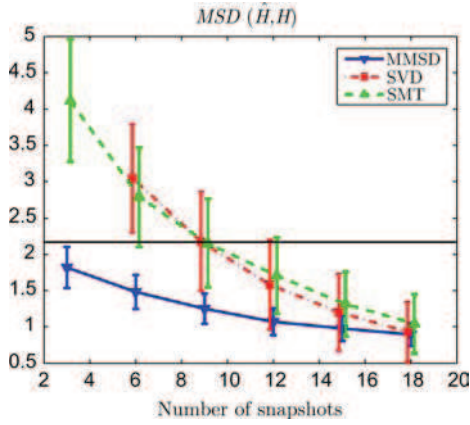


Fig. 5. Mean-square distance between true and estimated subspaces versus  $K$ .  $N = 20$ ,  $p = 5$ , SNR = 3 dB, and  $\theta_{\max} = 75^\circ$ .

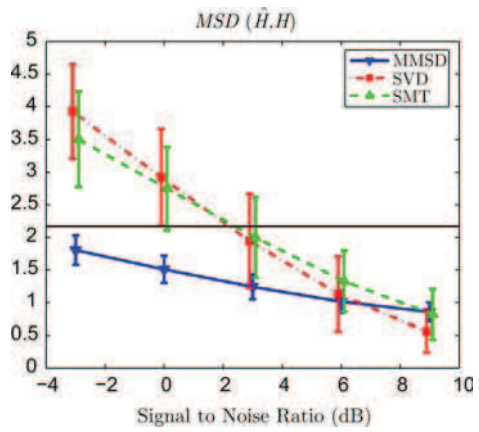


Fig. 8. Mean-square distance between true and estimated subspaces versus SNR.  $N = 20$ ,  $p = 5$ ,  $K = 10$ , and  $\theta_{\max} = 60^\circ$ .

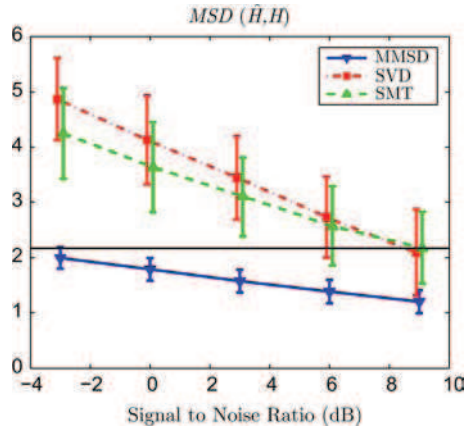


Fig. 6. Mean-square distance between true and estimated subspaces versus SNR.  $N = 20$ ,  $p = 5$ ,  $K = 5$ , and  $\theta_{\max} = 60^\circ$ .

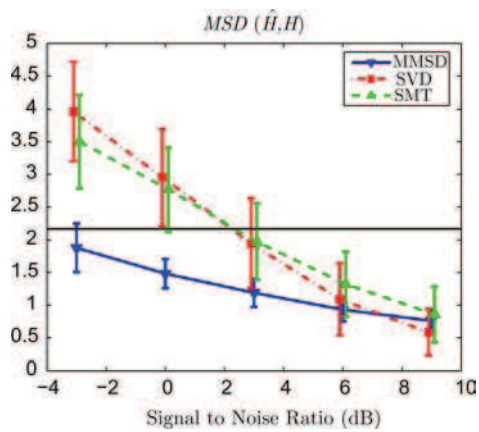


Fig. 9. Mean-square distance between true and estimated subspaces versus SNR.  $N = 20$ ,  $p = 5$ ,  $K = 10$ , and  $\theta_{\max} = 75^\circ$ .

of the soil characteristics, and thus numerous studies have focused on information retrieval from multi-band data; see, e.g., [26]–[29]. So far, a widely accepted model is that the image can be linearly decomposed as a combination of a few components, referred to as the endmembers [30]. One critical issue is thus to identify the subspace where the data lies together with the coordinates in this subspace, which provide the respective abundances, i.e., the proportion of the soil components. This

can be achieved by well-known and computationally efficient techniques such as principal component analysis (PCA), a primordial asset to using the linear (or subspace) model. However, it may be argued that the linear model does not fully account for all physical phenomenon that give rise to the image, e.g., the possibly non-linear mixing of the components. In order to obtain a finer image analysis, non-linear models can be investigated [30] but generally at the price of a higher computational



complexity. Furthermore, in most cases non-linear effects are not that important and an interesting alternative is to continue to resort to a linear model but at a local level (i.e., within a few pixels) rather than at the full image level. Doing so, one can characterize the data locally and track the evolution of the local subspaces in order to assess the degree of non-linearity. The subspace estimation scheme developed above can fulfill this task and it is now tested against real hyperspectral data, acquired by the NASA spectro-imager AVIRIS over Moffett Field, CA, in 1997. More precisely, we consider a  $50 \times 50$  sub-image, which contains partly a lake (upper part of the sub-image) and partly a coastal area (lower part of the sub-image) composed of soil and vegetation, see [31] for a more detailed description. The data is collected in  $N = 183$  spectral bands and we have thus a total of  $L = 2500$  pixels. Under the linear mixing model and in the absence of noise, the data matrix  $\mathbf{Y} = [\mathbf{y}_1 \cdots \mathbf{y}_L]$ , where  $\mathbf{y}_\ell \in \mathbb{R}^N$  stands for the  $\ell$ th pixel, can be written as  $\mathbf{Y} = \mathbf{M}\mathbf{A}$  where  $\mathbf{M} = [\mathbf{m}_1 \cdots \mathbf{m}_R]$  and  $\mathbf{m}_r$ ,  $r = 1, \dots, R$  denotes the set of endmembers, i.e., the spectral signatures which best describe the soil components. In [31], it was shown that a value  $R = 3$  was sufficient to obtain an accurate description of the data. The columns  $\mathbf{a}_\ell = [a_{\ell,1} \cdots a_{\ell,R}]^T$  of the matrix  $\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_L]$  are the so-called abundances: They satisfy the positivity constraint  $a_{\ell,r} \geq 0$  and the sum-to-one property, i.e.,  $\mathbf{a}_\ell^T \mathbf{1}_R = 1$  where  $\mathbf{1}_R$  is the  $R$ -length vector whose elements are all equal to 1. In other words, the matrix  $\mathbf{A}$  satisfies the constraint  $\mathbf{A}^T \mathbf{1}_R = \mathbf{1}_L$ . The pixels  $\mathbf{y}_\ell$  thus belong to a simplex whose vertices are the  $R$  endmembers  $\mathbf{m}_r$  [31]. Let  $\boldsymbol{\mu} = L^{-1} \sum_{\ell=1}^L \mathbf{y}_\ell$  denote the mean value of the pixels. Then, the centered data matrix  $\mathbf{X} = \mathbf{Y} - \boldsymbol{\mu} \mathbf{1}_L^T$  belongs to a  $p$ -dimensional subspace (with  $p = R - 1$ ), which can be estimated by a number of techniques, including PCA [31].

Usually, PCA is performed on the whole image, which makes sense if the linear mixing model is in force for all pixels. Herein, we are interested in assessing the validity of this model at the pixel level. More precisely, the PCA on the whole image will provide us with the ‘‘average’’ subspace: the pixels  $\mathbf{y}_\ell$  are then unitarily transformed ( $\mathbf{y}_\ell \leftarrow \mathbf{Q}\mathbf{y}_\ell$ ) such that  $\tilde{\mathbf{H}} \leftarrow \mathbf{Q}\tilde{\mathbf{H}} = [\mathbf{I}_p \ \mathbf{0}]^T$ , and we are interested in the distance between  $\tilde{\mathbf{H}}$  and the subspace spanned by a pixel  $\mathbf{y}_\ell$  and its few nearest pixels. If this distance is very small, then it is likely that the linear model described by  $\tilde{\mathbf{H}}$  is rather accurate. On the other hand, if the distance is not negligible, it may be that  $\tilde{\mathbf{H}}$  does not describe accurately the scene around pixel  $\ell$  or that some non-linear mixing effects might occur there. Therefore, subspace estimation at the pixel level together with distance to  $\tilde{\mathbf{H}}$  evaluation enables one to gain insight into the understanding of the mixing process. This is the approach we take here and our MMSD estimator is used towards this end. To be more specific, for each pixel  $\mathbf{y}_\ell$  we use the latter and its three nearest neighbors (hence  $K = 4$ ) to obtain the MMSD estimator of the local subspace. The mean square distance between  $\mathcal{R}(\hat{\mathbf{H}}_\ell)$  and  $\mathcal{R}(\tilde{\mathbf{H}})$ ,  $\text{MSD}(\hat{\mathbf{H}}_\ell, \tilde{\mathbf{H}})$  is then determined to evaluate how close are the local subspace and the global subspace. The results are shown in Fig. 10<sup>5</sup>: For comparison purposes, we display in

<sup>5</sup>Application to another image and results with a different value of  $K$  can be found in [24].

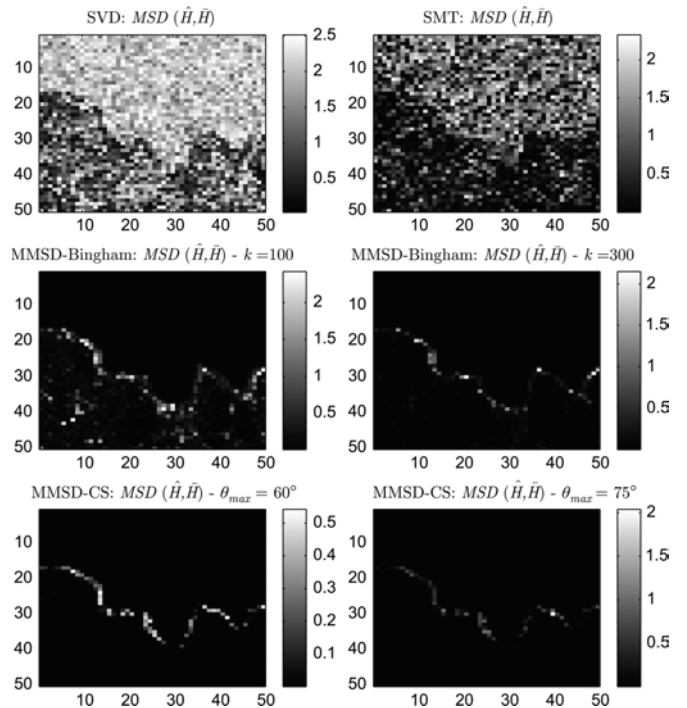


Fig. 10. Moffett image.  $\text{MSD}(\hat{\mathbf{H}}_\ell, \tilde{\mathbf{H}})$ .  $N = 183$ ,  $p = 2$ ,  $K = 4$ .

this figure the result obtained with the SVD, the SMT and the method of [17], which assumes a Bingham prior distribution for  $\mathbf{H}$ . Fig. 10 shows that a local SVD or SMT would predict rather large differences between the local subspaces and  $\tilde{\mathbf{H}}$ , especially for pixels in the lake area. However, it cannot be concluded that  $\tilde{\mathbf{H}}$  does not apply for most of the image since, with  $K = 4$ , the subspace estimated by the SVD may not be very accurate. In contrast, the Bayesian CS-based MMSD estimator shows that  $\tilde{\mathbf{H}}$  is a rather accurate subspace for the whole image (especially on the lake), except for the pixels along the transition between lake and coastal area. This seems logical as non-linear mixing effects are more likely to occur along the shore, while the linear model is likely to apply well elsewhere. Therefore, the MMSD estimator is able to reveal the zones of the image where departure from the linear model might occur. Finally, we note that it is not intuitive to set of value for  $\kappa$ : The values  $\kappa = 100$  and  $\kappa = 300$  do not have a real meaning and lead to different interpretations of the image. It is much easier to set a value for  $\theta_{\max}$ , a significant advantage of the CS-based model compared to the method of [17]. However, the latter is computationally less intensive. As a final comment, we would like to point out that the computational complexity of the present MMSD-CS method could be prohibitive in large dimensional problems ( $N$  large), for which more computationally efficient algorithms, such as the sparse matrix transform of [8], should be favored.

## V. CONCLUSION

In this paper, we considered the problem of subspace estimation from a possibly very limited number of snapshots under the assumption that some prior knowledge about the subspace

is available. A Bayesian statistical model was formulated to account for this situation, based on the CS decomposition of the semi-orthogonal matrix  $\mathbf{H}$  whose columns span the subspace of interest. This model was shown to rely on rather mild assumptions and, moreover, these assumptions involve meaningful and intuitively appealing quantities, namely the angles between the prior subspace  $\tilde{\mathbf{H}}$  and the true subspace  $\mathbf{H}$ . The minimum mean-square distance estimator was implemented through a Gibbs sampling scheme. It was shown to provide accurate estimates, in particular in the low SNR or low sample support regimes. The estimator was also successfully applied to real hyperspectral data, demonstrating its ability to reveal the limits of linear mixing models.

## REFERENCES

- [1] L. L. Scharf, *Statistical Signal Processing: Detection, Estimation and Time Series Analysis*. Reading, MA: Addison-Wesley, 1991.
- [2] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1993.
- [3] R. Kumaresan and D. Tufts, "Estimating the parameters of exponentially damped sinusoids and pole-zero modeling in noise," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 30, no. 6, pp. 833–840, Dec. 1982.
- [4] R. Kumaresan and D. Tufts, "Estimating the angles of arrival of multiple plane waves," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 19, no. 1, pp. 134–139, Jan. 1983.
- [5] P. Stoica and A. Nehorai, "MUSIC, maximum likelihood and Cramér–Rao bound," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 37, no. 5, pp. 720–741, May 1989.
- [6] O. Ledoit and M. Wolf, "A well-conditioned estimator for large-dimensional covariance matrices," *J. Multivar. Anal.*, vol. 88, no. 2, pp. 365–411, Feb. 2004.
- [7] T. L. Marzetta, G. H. Tucci, and S. H. Simon, "A random matrix theoretic approach to handling singular covariance estimates," *IEEE Trans. Inf. Theory*, vol. 57, no. 9, pp. 6256–6271, Sep. 2011.
- [8] G. Cao, L. R. Bachega, and C. A. Bouman, "The sparse matrix transform for covariance estimation and analysis of high dimensional signals," *IEEE Trans. Image Process.*, vol. 20, no. 3, pp. 625–640, Mar. 2011.
- [9] X. Mestre, "Improved estimation of eigenvalues and eigenvectors of covariance matrices using their sample estimates," *IEEE Trans. Inf. Theory*, vol. 54, no. 11, pp. 5113–5129, Nov. 2008.
- [10] J. Thomas, L. Scharf, and D. Tufts, "The probability of a subspace swap in the SVD," *IEEE Trans. Signal Process.*, vol. 43, no. 3, pp. 730–736, Mar. 1995.
- [11] M. Hawkes, A. Nehorai, and P. Stoica, "Performance breakdown of subspace-based methods: prediction and cure," *Proc. Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, pp. 4005–4008, May 2001.
- [12] Z. Bai and J. W. Silverstein, *Spectral Analysis of Large Dimensional Random Matrices*, ser. Springer Series in Statistics, 2nd ed. New York: Springer-Verlag, 2010.
- [13] D. Paul, "Asymptotics of sample eigenstructure for a large dimensional spiked covariance model," *Stat. Sinica*, vol. 17, no. 4, pp. 1617–1642, Oct. 2007.
- [14] F. Benaych-Georges and R. R. Nadakuditi, "The singular values and vectors of low rank perturbations of large rectangular random matrices," 2011 [Online]. Available: <http://arxiv.org/abs/1103.2221>
- [15] F. Benaych-Georges and R. R. Nadakuditi, "The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices," *Adv. Math.*, vol. 211, no. 1, pp. 494–521, May 2011.
- [16] A. Srivastava, "A Bayesian approach to geometric subspace estimation," *IEEE Trans. Signal Process.*, vol. 48, no. 5, pp. 1390–1400, May 2000.
- [17] O. Besson, N. Dobigeon, and J.-Y. Tournet, "Minimum mean square distance estimation of a subspace," *IEEE Trans. Signal Process.*, vol. 59, no. 12, pp. 5709–5720, Dec. 2011.
- [18] A. Edelman, T. Arias, and S. Smith, "The geometry of algorithms with orthogonality constraints," *SIAM J. Matrix Anal. Appl.*, vol. 20, no. 2, pp. 303–353, 1998.
- [19] G. Golub and C. V. Loan, *Matrix Computations*, 3rd ed. Baltimore, MD: The John Hopkins Univ. Press, 1996.
- [20] K. V. Mardia and P. E. Jupp, *Directional Statistics*. New York: Wiley, 1999.
- [21] C. P. Robert and G. Casella, *Monte Carlo Statistical Methods*, 2nd ed. New York: Springer-Verlag, 2004.
- [22] C. P. Robert, *The Bayesian Choice—From Decision-Theoretic Foundations to Computational Implementation*. New York: Springer-Verlag, 2007.
- [23] P. D. Hoff, "Simulation of the matrix Bingham–Von Mises–Fisher distribution, with applications to multivariate and relational data," *J. Comput. Graph. Stat.*, vol. 18, no. 2, pp. 438–456, Jun. 2009.
- [24] O. Besson, N. Dobigeon, and J.-Y. Tournet, "CS decomposition based Bayesian subspace estimation," IIRIT/ENSEEIH, Toulouse, France, 2012.
- [25] C.-I Chang, *Hyperspectral Imaging: Techniques for Spectral Detection and Classification*. New York: Kluwer, 2003.
- [26] D. Manolakis, C. Siracusa, and G. Shaw, "Hyperspectral subpixel target detection using the linear mixing model," *IEEE Trans. Geosci. Remote Sens.*, vol. 39, no. 7, pp. 1392–409, Jul. 2001.
- [27] M. Lewis, V. Jooste, and A. A. de Gasparis, "Discrimination of arid vegetation with airborne multispectral scanner hyperspectral imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 39, no. 7, pp. 1471–1479, Jul. 2001.
- [28] B. Datt, T. R. McVicar, T. G. V. Niel, D. L. B. Jupp, and J. S. Pearlman, "Preprocessing EO-1 hyperion hyperspectral data to support the application of agricultural indexes," *IEEE Trans. Geosci. Remote Sens.*, vol. 41, no. 6, pp. 1246–1259, Jun. 2003.
- [29] J. Plaza, R. Pérez, A. Plaza, P. Martínez, and D. Valencia, "Mapping oil spills on sea water using spectral mixture analysis of hyperspectral image data," in *Chemical and Biological Standoff Detection III*, J. O. Jensen and J.-M. Thériault, Eds. Bellingham, WA: SPIE, 2005, vol. 5995, pp. 79–86.
- [30] N. Keshava and J. Mustard, "Spectral unmixing," *IEEE Signal Process. Mag.*, vol. 19, no. 1, pp. 44–57, Jan. 2002.
- [31] N. Dobigeon, S. Moussaoui, M. Coulon, J.-Y. Tournet, and A. O. Hero, "Joint Bayesian endmember extraction and linear unmixing for hyperspectral imagery," *IEEE Trans. Signal Process.*, vol. 57, no. 11, pp. 4355–4368, Nov. 2009.



**Olivier Besson** (S'90–M'93–SM'04) received the Ph.D. degree in signal processing and the Habilitation à Diriger des Recherches from INP Toulouse, France, in 1992 and 1998, respectively.

He is currently a Professor with the Department of Electronics, Optronics and Signal of the Institut Supérieur de l'Aéronautique et de l'Espace (ISAE), Toulouse, France. His research interests are in the area of robust adaptive array processing, mainly for radar applications.

Dr. Besson is a former Associate Editor of the IEEE TRANS. SIGNAL PROCESS. and the IEEE SIGNAL PROCESSING LETTERS. He is a member of the Sensor Array and Multichannel Technical Committee (SAM TC) of the IEEE Signal Processing Society.



**Nicolas Dobigeon** (S'05–M'08) was born in Angoulême, France, in 1981. He received the Eng. degree in electrical engineering from ENSEEIHT, Toulouse, France, and the M.Sc. degree in signal processing from the National Polytechnic Institute of Toulouse, France, both in 2004, and the Ph.D. degree in signal processing from the National Polytechnic Institute of Toulouse, France, in 2007.

From 2007 to 2008, he was a Postdoctoral Research Associate at the Department of Electrical Engineering and Computer Science, University of Michigan. Since 2008, he has been an Assistant Professor with the National Polytechnic Institute of Toulouse (ENSEEIH—University of Toulouse), France, within the Signal and Communication Group of the IIRIT Laboratory. His research interests are centered around statistical signal and image processing, with a particular interest in Bayesian inference and Markov chain Monte Carlo (MCMC) methods.



**Jean-Yves Tourneret** (SM'08) received the Ingénieur degree in electrical engineering from the Ecole Nationale Supérieure d'Electronique, d'Electrotechnique, d'Informatique et d'Hydraulique, Toulouse (ENSEEIH), France, in 1989 and the Ph.D. degree from the National Polytechnic Institute, Toulouse, France, in 1992.

He is currently a Professor in the University of Toulouse (ENSEEIH), France, and a member of the IRIT laboratory (UMR 5505 of the CNRS). His research activities are centered around statistical

signal processing, with a particular interest to Bayesian and Markov chain Monte Carlo methods.

Dr. Tourneret has been involved in the organization of several conferences, including the European Conference on Signal Processing (EUSIPCO) 2002 (as the program chair), the International Conference on Acoustics, Speech and Signal Processing (ICASSP) 2006 (in charge of plenaries) and the Statistical Signal Processing Workshop (SSP) 2012 (for international liaisons). He has been a member of different technical committees, including the Signal Processing Theory and Methods (SPTM) Committee of the IEEE Signal Processing Society from 2001 to 2007 and from 2010 to present. He served as an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING from 2008 to 2011.