

# Minimum mean square distance estimation of a subspace

Olivier Besson\*, Nicolas Dobigeon† and Jean-Yves Tournet†

Technical report - 14 June 2011

## Abstract

We consider the problem of subspace estimation in a Bayesian setting. Since we are operating in the Grassmann manifold, the usual approach which consists of minimizing the mean square error (MSE) between the true subspace  $\mathbf{U}$  and its estimate  $\hat{\mathbf{U}}$  may not be adequate as the MSE is not the natural metric in the Grassmann manifold  $G_{N,p}$ , i.e., the set of  $p$ -dimensional subspaces in  $\mathbb{R}^N$ . As an alternative, we propose to carry out subspace estimation by minimizing the mean square distance between  $\mathbf{U}$  and its estimate, where the considered distance is a natural metric in the Grassmann manifold, viz. the distance between the projection matrices. We show that the resulting estimator is no longer the posterior mean of  $\mathbf{U}$  but entails computing the principal eigenvectors of the posterior mean of  $\mathbf{U}\mathbf{U}^T$ . Derivation of the minimum mean square distance (MMSD) estimator is carried out in a few illustrative examples including a linear Gaussian model for the data and Bingham or von Mises Fisher prior distributions for  $\mathbf{U}$ . In all scenarios, posterior distributions are derived and the MMSD estimator is obtained either analytically or implemented via a Markov chain Monte Carlo simulation method. The method is shown to provide accurate estimates even when the number of samples is lower than the dimension of  $\mathbf{U}$ . An application to hyperspectral imagery is finally investigated.

---

\*O. Besson is with the University of Toulouse, ISAE, Department Electronics Optonics Signal, Toulouse (e-mail: olivier.besson@isae.fr). The work of O. Besson was partly supported by DGA-MRIS under grant no. 2009.60.033.00.470.75.01.

†N. Dobigeon and J.-Y. Tournet are with the University of Toulouse, IRIT/ENSEEIH, Signal and Communications Group, Toulouse, France (e-mail: nicolas.dobigeon@enseeiht.fr, jean-yves.tournet@enseeiht.fr).

# 1 Problem statement

In many signal processing applications, the signals of interest do not span the entire observation space and a relevant and frequently used assumption is that they evolve in a low-dimensional subspace [1]. Subspace modeling is accurate when the signals consist of a linear combination of  $p$  modes in an  $N$ -dimensional space, and constitute a good approximation for example when the signal covariance matrix is close to rank-deficient. As a consequence, subspace estimation plays a central role in recovering these signals with maximum accuracy. Frequency estimation of complex exponential signals and directions of arrival (DoA) of signals using an array of sensors are examples of well-known area where obtaining accurate subspace estimates constitutes a crucial task. In the latter case, the data matrix can be written as  $\mathbf{Y} = \mathbf{A}(\boldsymbol{\theta})\mathbf{S} + \mathbf{N}$  where the columns of the  $N \times p$  matrix  $\mathbf{A}(\boldsymbol{\theta})$  contain the  $p$  steering vectors of the sources of interest, and  $\boldsymbol{\theta}$  corresponds to the vector of their DoA. In a noise-free environment,  $\mathbf{Y}$  is rank-deficient and its column space is spanned by the steering vectors, from which DoA estimation can be performed. Therefore, a key issue is to first identify the dominant subspace of  $\mathbf{Y}$  in order then to infer the DoA. This is the essence of most subspace-based methods, such as MUSIC [2], ESPRIT [3] or Min-Norm [4,5]. These methods have gained much popularity because they represent computationally interesting alternatives to the maximum likelihood (ML) estimator [6]. Moreover, they yield equivalent performance in the asymptotic regime, viz large number of snapshots or high signal to noise ratio (SNR) [7,8]. An ubiquitous tool to compute the dominant subspace is the singular value decomposition (SVD) of the data matrix or, equivalently, the eigenvalue decomposition of the sample covariance matrix. Observe also that the SVD emerges naturally as the maximum likelihood estimator of the range space  $\mathcal{R}(\mathbf{U})$  of  $\mathbf{U}$  in the classical model  $\mathbf{Y} = \mathbf{U}\mathbf{S} + \mathbf{N}$ , where  $\mathbf{Y}$  stands for the  $N \times K$  observation matrix,  $\mathbf{U}$  is the (deterministic)  $N \times p$  matrix, with  $p < N$ , whose columns span the  $p$ -dimensional subspace of interest,  $\mathbf{S}$  is the  $p \times K$  (deterministic) waveform matrix and  $\mathbf{N}$  is the additive noise. The latter model is relevant e.g., in hyperspectral imagery, see Section 5 for further discussion.

When the SNR is high or the number of snapshots is large, the SVD yields quasi-optimal estimates. However, the SVD can incur some performance loss in two main cases, namely when the SNR is very low and thereof the probability of a subspace swap or subspace leakage is high [9–12]. A second case occurs when the number of samples  $K$  is small, typically of the order of  $N$  or, even less, of the order of  $p$ . In order to overcome these limitations, especially the case of low sample support, some interesting alternatives have been proposed, such as [13], where accurate estimates of the eigenvalues and eigenvectors are derived, based on random matrix theory. However, the sample support considered there is still higher than what we consider herein. More precisely, we are mostly interested in the case where the number of snapshots is very small and may be less than the subspace dimension  $p$ : in this case,  $\mathbf{Y}$  is at most of rank  $K$  and information is lacking about how to complement  $\mathcal{R}(\mathbf{Y})$  in order to estimate  $\mathcal{R}(\mathbf{U})$ . This case along with the low SNR case are those of most interest to us and, hence, we need to turn to a different methodology.

In such situations, a Bayesian approach might be helpful as it enables one to assist estimation by providing some statistical information about  $\mathbf{U}$ . We investigate such an approach herein and assign to the unknown matrix  $\mathbf{U}$  an appropriate prior distribution, taking into account the specific structure of  $\mathbf{U}$ . The paper is organized as follows. In section 2, we propose an approach based on minimizing a natural distance on the Grassmann manifold, which yields a new estimator of  $\mathbf{U}$ . The theory is illustrated in section 3 where the new estimator is derived for some specific examples. In section 4 its performance is assessed through numerical simulations, and compared with conventional approaches. Section 5 studies an application to the analysis of interactions between pure materials contained in hyperspectral images.

## 2 Minimum mean square distance estimation

As explained above, we adopt a Bayesian framework for estimation of the unknown matrix  $\mathbf{U}$  and consider an alternative to the conventional minimum mean square error (MMSE) estimator, which is usually regarded as the chief systematic approach under the Bayesian umbrella [14]. Let us consider that we wish to estimate the range space of  $\mathbf{U}$  from the joint distribution  $p(\mathbf{Y}, \mathbf{U})$  where  $\mathbf{Y}$  stands for the available data matrix. Usually, one is not interested in  $\mathbf{U}$  *per se* but rather in its range space  $\mathcal{R}(\mathbf{U})$ , and thus we are operating in the Grassmann manifold  $G_{N,p}$ , i.e., the set of  $p$ -dimensional subspaces in  $\mathbb{R}^N$  [15]. It is thus natural to wonder whether the MMSE estimator is suitable in  $G_{N,p}$ . The MMSE estimator  $\hat{\boldsymbol{\theta}}$  of a vector  $\boldsymbol{\theta}$  minimizes the average squared Euclidean distance between  $\hat{\boldsymbol{\theta}}$  and  $\boldsymbol{\theta}$ , i.e.,  $\mathbb{E} \left\{ \left\| \hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \right\|_2^2 \right\}$ . Despite the fact that this distance is natural in an Euclidean space, it may not be the more natural metric in  $G_{N,p}$ . In fact, the natural distance between two subspaces  $\mathcal{R}(\mathbf{U}_1)$  and  $\mathcal{R}(\mathbf{U}_2)$  is given by  $(\sum_{k=1}^p \theta_k^2)^{1/2}$  [15] where  $\theta_k$  are the principal angles between these subspaces, which can be obtained by SVD of  $\mathbf{U}_2^T \mathbf{U}_1$  where  $\mathbf{U}_1$  and  $\mathbf{U}_2$  denote orthonormal bases for these subspaces [16]. The SVD of  $\mathbf{U}_2^T \mathbf{U}_1$  is defined as  $\mathbf{U}_2^T \mathbf{U}_1 = \mathbf{X} \text{diag}(\cos \theta_1, \dots, \cos \theta_p) \mathbf{Z}^T$ , where  $\mathbf{X}$  and  $\mathbf{Z}$  are two  $p \times p$  unitary matrices. Therefore, it seems more adequate, rather than minimizing  $\left\| \hat{\mathbf{U}} - \mathbf{U} \right\|_F^2$  as the MMSE estimator does, to minimize the natural distance between the subspaces spanned by  $\hat{\mathbf{U}}$  and  $\mathbf{U}$ . Although this is the most intuitively appealing method, it faces the drawback that the cosines of the angles and not the angles themselves emerge naturally from the SVD. Therefore, for the sake of practicality, we consider minimizing the sum of the squared sine of the angles between  $\hat{\mathbf{U}}$  and  $\mathbf{U}$ . As argued in [15, 16], this cost function is natural in the Grassmann manifold since it corresponds to the Frobenius norm of the difference between the projection matrices on the two subspaces, viz  $\sum_{k=1}^p \sin^2 \theta_k = \left\| \hat{\mathbf{U}} \hat{\mathbf{U}}^T - \mathbf{U} \mathbf{U}^T \right\|_F^2 \triangleq d^2(\hat{\mathbf{U}}, \mathbf{U})$ . It should be mentioned that our approach follows along the same principles as in [17] where a Bayesian framework is proposed for subspace estimation, and where the author considers minimizing  $d(\hat{\mathbf{U}}, \mathbf{U})$ . Hence the theory presented in this section is similar to that of [17], only the parameterization of the problem in [17] being slightly different from ours. The main differences compared to [17] lie in the prior distributions and signal models used, as well as in the implementation of the MMSD estimator, see the next section for details.

Given that  $d^2(\hat{\mathbf{U}}, \mathbf{U}) = 2(p - \text{Tr} \{ \hat{\mathbf{U}}^T \mathbf{U} \mathbf{U}^T \hat{\mathbf{U}} \})$ , we define the minimum mean-square distance (MMSD) estimator of  $\mathbf{U}$  as

$$\hat{\mathbf{U}}_{\text{mmsd}} = \arg \max_{\hat{\mathbf{U}}} \mathbb{E} \left\{ \text{Tr} \left\{ \hat{\mathbf{U}}^T \mathbf{U} \mathbf{U}^T \hat{\mathbf{U}} \right\} \right\}. \quad (1)$$

Since

$$\mathbb{E} \left\{ \text{Tr} \left\{ \hat{\mathbf{U}}^T \mathbf{U} \mathbf{U}^T \hat{\mathbf{U}} \right\} \right\} = \int \left[ \int \text{Tr} \left\{ \hat{\mathbf{U}}^T \mathbf{U} \mathbf{U}^T \hat{\mathbf{U}} \right\} p(\mathbf{U} | \mathbf{Y}) d\mathbf{U} \right] p(\mathbf{Y}) d\mathbf{Y} \quad (2)$$

it follows that

$$\begin{aligned} \hat{\mathbf{U}}_{\text{mmsd}} &= \arg \max_{\hat{\mathbf{U}}} \int \text{Tr} \left\{ \hat{\mathbf{U}}^T \mathbf{U} \mathbf{U}^T \hat{\mathbf{U}} \right\} p(\mathbf{U} | \mathbf{Y}) d\mathbf{U} \\ &= \arg \max_{\hat{\mathbf{U}}} \text{Tr} \left\{ \hat{\mathbf{U}}^T \left[ \int \mathbf{U} \mathbf{U}^T p(\mathbf{U} | \mathbf{Y}) d\mathbf{U} \right] \hat{\mathbf{U}} \right\}. \end{aligned} \quad (3)$$

Therefore, the MMSD estimate of the subspace spanned by  $\mathbf{U}$  is given by the  $p$  largest eigenvectors of the matrix  $\int \mathbf{U} \mathbf{U}^T p(\mathbf{U} | \mathbf{Y}) d\mathbf{U}$ , which we denote as

$$\hat{\mathbf{U}}_{\text{mmsd}} = \mathcal{P}_p \left\{ \int \mathbf{U} \mathbf{U}^T p(\mathbf{U} | \mathbf{Y}) d\mathbf{U} \right\}. \quad (4)$$

In other words, MMSD estimation amounts to find the best rank- $p$  approximation to the posterior mean of the projection matrix  $\mathbf{U}\mathbf{U}^T$  on  $\mathcal{R}(\mathbf{U})$ . For notational convenience, let us denote  $\mathbf{M}(\mathbf{Y}) = \int \mathbf{U}\mathbf{U}^T p(\mathbf{U}|\mathbf{Y}) d\mathbf{U}$ . Except for a few cases where this matrix can be derived in closed-form (an example is given in the next section), there usually does not exist any analytical expression for  $\mathbf{M}(\mathbf{Y})$ . In such situation, an efficient way to approximate the matrix  $\mathbf{M}(\mathbf{Y})$  is to use a Markov chain Monte Carlo (MCMC) simulation method whose goal is to generate random matrices  $\mathbf{U}$  drawn from the posterior distribution  $p(\mathbf{U}|\mathbf{Y})$ , and to approximate the integral in (4) by a finite sum. This aspect will be further elaborated in the next section. Let  $\mathbf{M}(\mathbf{Y}) = \mathbf{U}_M(\mathbf{Y})\mathbf{L}_M(\mathbf{Y})\mathbf{U}_M^T(\mathbf{Y})$  denote the eigenvalue decomposition of  $\mathbf{M}(\mathbf{Y})$  with  $\mathbf{L}_M(\mathbf{Y}) = \text{diag}(\ell_1(\mathbf{Y}), \ell_2(\mathbf{Y}), \dots, \ell_N(\mathbf{Y}))$  and  $\ell_1(\mathbf{Y}) \geq \ell_2(\mathbf{Y}) \geq \dots \geq \ell_N(\mathbf{Y})$ . Then the average distance between  $\hat{\mathbf{U}}_{\text{mmsd}}$  and  $\mathbf{U}$  is given by

$$\begin{aligned} \mathbb{E} \left\{ d^2 \left( \hat{\mathbf{U}}_{\text{mmsd}}, \mathbf{U} \right) \right\} &= 2p - 2 \int \left[ \int \text{Tr} \left\{ \hat{\mathbf{U}}_{\text{mmsd}}^T \mathbf{U}\mathbf{U}^T \hat{\mathbf{U}}_{\text{mmsd}} \right\} p(\mathbf{U}|\mathbf{Y}) d\mathbf{U} \right] p(\mathbf{Y}) d\mathbf{Y} \\ &= 2p - 2 \int \text{Tr} \left\{ \hat{\mathbf{U}}_{\text{mmsd}}^T \mathbf{M}(\mathbf{Y}) \hat{\mathbf{U}}_{\text{mmsd}} \right\} p(\mathbf{Y}) d\mathbf{Y} \\ &= 2p - 2 \sum_{k=1}^p \int \ell_k(\mathbf{Y}) p(\mathbf{Y}) d\mathbf{Y}. \end{aligned} \quad (5)$$

The latter expression constitutes a lower bound on  $\mathbb{E} \left\{ d^2 \left( \hat{\mathbf{U}}, \mathbf{U} \right) \right\}$  and is referred to as the Hilbert-Schmidt bound in [17, 18]. As indicated in these references, and similarly to  $\mathbf{M}(\mathbf{Y})$ , this lower bound may be difficult to obtain analytically.

The MMSD approach can be extended to the mixed case where, in addition to  $\mathbf{U}$ , a parameter vector  $\boldsymbol{\theta}$  which can take arbitrary values in  $\mathbb{R}^q$  needs to be estimated jointly with  $\mathbf{U}$ . Under such circumstances, one can estimate  $\mathbf{U}$  and  $\boldsymbol{\theta}$  as

$$\left( \hat{\mathbf{U}}_{\text{mmsd}}, \hat{\boldsymbol{\theta}}_{\text{mmsd}} \right) = \arg \min_{\hat{\mathbf{U}}, \hat{\boldsymbol{\theta}}} \mathbb{E} \left\{ -\text{Tr} \left\{ \hat{\mathbf{U}}^T \mathbf{U}\mathbf{U}^T \hat{\mathbf{U}} \right\} + \left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \right)^T \left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \right) \right\}. \quad (6)$$

Doing so, the MMSD estimator of  $\mathbf{U}$  is still be given by (4) while the MMSD and MMSE estimators of  $\boldsymbol{\theta}$  coincide.

*Remark 1.* The MMSD approach differs from an MMSE approach which would entail calculating the posterior mean of  $\mathbf{U}$ , viz  $\int \mathbf{U} p(\mathbf{U}|\mathbf{Y}) d\mathbf{U}$ . Note that the latter may not be meaningful, in particular when the posterior distribution  $p(\mathbf{U}|\mathbf{Y})$  depends on  $\mathbf{U}$  only through  $\mathbf{U}\mathbf{U}^T$ , see next section for an example. In such a case, post-multiplication of  $\mathbf{U}$  by any  $p \times p$  unitary matrix  $\mathbf{Q}$  yields the same value of  $p(\mathbf{U}|\mathbf{Y})$ . Therefore averaging  $\mathbf{U}$  over  $p(\mathbf{U}|\mathbf{Y})$  does not make sense while computing (4) is relevant. On the other hand, if  $p(\mathbf{U}|\mathbf{Y})$  depends on  $\mathbf{U}$  directly, then computing the posterior mean of  $\mathbf{U}$  can be investigated: an example where this situation occurs will be presented in the next section. As a final comment, observe that  $\int \mathbf{U} p(\mathbf{U}|\mathbf{Y}) d\mathbf{U}$  is not necessarily unitary but its range space can be used to estimate  $\mathcal{R}(\mathbf{U})$ .

*Remark 2.* We open a parenthesis here regarding the framework of this paper. First of all, we do not deal herein with sequential estimation in a Bayesian framework as can be the case e.g., in [17, 19]. In contrast, we consider batch algorithms where a given number of snapshots is used to produce an estimate of  $\mathbf{U}$ . A second important notice is that the MMSD estimator of this paper originally stems from the minimization of a certain distance over the Grassmann manifold or the Stiefel manifold (the set of  $N \times p$  matrices  $\mathbf{U}$  such that  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ ). Optimization over these manifolds has received considerable attention recently, see the excellent tutorial paper by Edelman *et al.* [15] as well as [20, 21], and [22–24] for signal processing applications. More generally, minimization over special manifolds (see e.g., [25] for optimization over the special

linear group) is attracting interest as it enables the underlying geometry of the problem to be taken into account. In the aforementioned references, the focus is on deriving iterative methods, such as steepest descent, conjugate gradient or Newton methods. In the present paper, the problem addressed is slightly different: we do not consider optimization with respect to an unknown matrix  $\mathbf{U}$ , rather we need to compute the average of  $\mathbf{U}\mathbf{U}^T$  over  $p(\mathbf{U}|\mathbf{Y})$  when  $\mathbf{U}$  is a random matrix on the Stiefel manifold.

### 3 Illustration examples

In this section we illustrate the previous theory on some examples, including the conventional linear Gaussian model (conditioned on  $\mathbf{U}$ ) and a model involving the eigenvalue decomposition of the data covariance matrix. As a first step, we address the issue of selecting prior distributions for  $\mathbf{U}$  and then move on to the derivation of the MMSD estimator.

#### 3.1 Prior distributions

A crucial step in any Bayesian estimation scheme consists of selecting the prior distribution for the variables to be estimated. We focus here on distributions on the Stiefel or Grassmann manifold, depending whether we consider the matrix  $\mathbf{U}$  itself or its range space. There exist only a few distributions on the Stiefel or Grassmann manifolds, the most widely accepted being the Bingham or von Mises Fisher (vMF) distributions [26, 27], which are given respectively by

$$p_{\text{B}}(\mathbf{U}) = \frac{1}{{}_1F_1\left(\frac{1}{2}p, \frac{1}{2}N; \mathbf{A}\right)} \text{etr}\{\mathbf{U}^T \mathbf{A} \mathbf{U}\} \quad (7)$$

$$p_{\text{vMF}}(\mathbf{U}) = \frac{1}{{}_0F_1\left(\frac{1}{2}N; \frac{1}{4}\mathbf{F}^T \mathbf{F}\right)} \text{etr}\{\mathbf{F}^T \mathbf{U}\} \quad (8)$$

where  $\text{etr}\{\cdot\}$  stands for the exponential of the trace of the matrix between braces,  $\mathbf{A}$  is an  $N \times N$  symmetric matrix,  $\mathbf{F}$  is an  $N \times p$  arbitrary matrix, and  ${}_0F_1(a; \mathbf{X})$ ,  ${}_1F_1(a, b; \mathbf{X})$  are hypergeometric functions of matrix arguments, see e.g., [27] for their definitions. The Bingham and the von Mises Fisher distributions have been proposed in various applications, including meteorology, biology, medicine, image analysis (see [26] and references therein), and recently for modeling of multipath communications channels [28, 29] or in shape analysis [30]. We will denote these distributions as  $\text{B}(\mathbf{A})$  and  $\text{vMF}(\mathbf{F})$ , respectively. Observe that the Bingham distribution depends on  $\mathbf{U}\mathbf{U}^T$  only, and can thus be viewed as a distribution on the Grassmann manifold [26, 27] while the vMF distribution depends on  $\mathbf{U}$  and is a distribution on the Stiefel manifold. In most applications, one is mostly interested in the projection matrix  $\mathbf{U}\mathbf{U}^T$  and therefore the Bingham distribution appears to be a more natural choice than the vMF distribution. In our case, in order to introduce some knowledge about  $\mathbf{U}$ , we assume that it is “close” to a given subspace spanned by the columns of an orthonormal matrix  $\bar{\mathbf{U}}$ , and hence we consider two possible prior distributions for  $\mathbf{U}$ , namely

$$\pi_{\text{B}}(\mathbf{U}) \propto \text{etr}\left\{\kappa \mathbf{U}^T \bar{\mathbf{U}} \bar{\mathbf{U}}^T \mathbf{U}\right\} \quad (9)$$

$$\pi_{\text{vMF}}(\mathbf{U}) \propto \text{etr}\left\{\kappa \mathbf{U}^T \bar{\mathbf{U}}\right\} \quad (10)$$

where  $\propto$  means “proportional to”. The matrix  $\bar{\mathbf{U}}$  reflects our knowledge about the subspace where the signals evolve. This matrix can be obtained from the data itself (see section 5.2 for an example about hyperspectral imagery) or from some available models. For instance, in radar applications, if the clutter subspace is to be estimated, there exists a number of relevant models, including the general covariance matrix model of [31], that can be used to obtain  $\bar{\mathbf{U}}$ .

The distribution in (9) is proportional to the sum of the squared cosine angles between  $\mathcal{R}(\mathbf{U})$  and  $\mathcal{R}(\bar{\mathbf{U}})$  while  $\pi_{\text{vMF}}(\mathbf{U})$  is proportional to the sum of the cosine angles between  $\mathcal{R}(\mathbf{U})$  and  $\mathcal{R}(\bar{\mathbf{U}})$ . Note that  $\kappa$  is a concentration parameter: the larger  $\kappa$  the more concentrated around  $\bar{\mathbf{U}}$  are the subspaces  $\mathbf{U}$ . The difference between the two distributions is the following. In the Bingham distribution only  $\mathcal{R}(\mathbf{U})$  and  $\mathcal{R}(\bar{\mathbf{U}})$  are close (at least for large values of  $\kappa$ ) since  $\pi_{\text{B}}(\mathbf{U})$  is invariant to post-multiplication of  $\mathbf{U}$  by any  $p \times p$  unitary matrix  $\mathbf{Q}$ . Hence  $\mathbf{U}$  is not necessarily close to  $\bar{\mathbf{U}}$ . In contrast, under the vMF prior distribution,  $\mathbf{U}$  and  $\bar{\mathbf{U}}$  are close. For illustration purposes, Figure 1 displays the average fraction of energy of  $\mathbf{U}$  in  $\mathcal{R}(\bar{\mathbf{U}})$  defined as

$$\text{AFE}(\mathbf{U}, \bar{\mathbf{U}}) = \mathbb{E} \left\{ \text{Tr} \left\{ \mathbf{U}^T \bar{\mathbf{U}} \bar{\mathbf{U}}^T \mathbf{U} \right\} / p \right\}. \quad (11)$$

As can be observed from these figures, both distributions allow the distance between  $\mathbf{U}$  and  $\bar{\mathbf{U}}$  to be set in a rather flexible way. Their AFE is shown to be identical for small values of the concentration parameter but, when  $\kappa$  increases, the AFE of the vMF distribution increases faster. Additionally, even if the AFE are close for small values of  $\kappa$ , the distributions of the

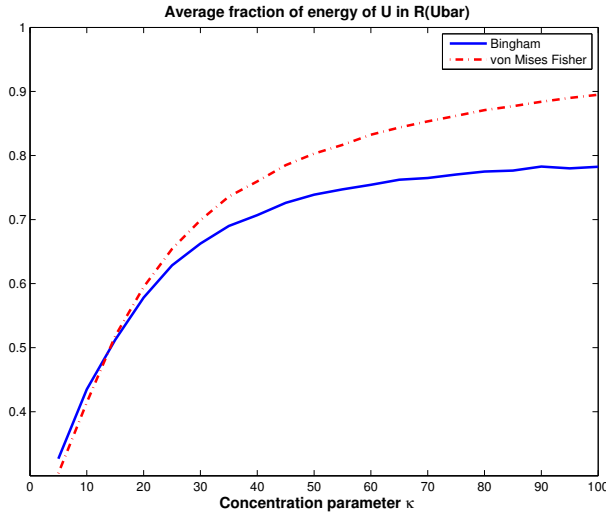


Figure 1: Average fraction of energy of  $\mathbf{U}$  in  $\mathcal{R}(\bar{\mathbf{U}})$  versus  $\kappa$ .  $N = 20$ ,  $p = 5$ .

angles between  $\mathcal{R}(\mathbf{U})$  and  $\mathcal{R}(\bar{\mathbf{U}})$  exhibit some differences, as shown in Figures 2 and 3 which display the probability density functions of these angles for  $\kappa = 20$ .

### 3.2 Linear model

In order to illustrate how the previous theory can be used in practice, we first consider a simple example, namely a linear Gaussian model (conditioned on  $\mathbf{U}$ ), i.e., we assume that the data follows the model  $\mathbf{Y} = \mathbf{U}\mathbf{S} + \mathbf{N}$  where the columns of  $\mathbf{N}$  are independent and identically distributed (i.i.d.) Gaussian vectors with zero-mean and (known) covariance matrix  $\sigma_n^2 \mathbf{I}$ . We assume that no knowledge about  $\mathbf{S}$  is available and hence its prior distribution is set to  $\pi(\mathbf{S}) \propto 1$ . Therefore, conditioned on  $\mathbf{U}$  we have

$$\begin{aligned} p(\mathbf{Y}|\mathbf{U}) &= \int p(\mathbf{Y}|\mathbf{U}, \mathbf{S}) \pi(\mathbf{S}) d\mathbf{S} \\ &\propto \int \text{etr} \left\{ -\frac{1}{2\sigma_n^2} (\mathbf{Y} - \mathbf{U}\mathbf{S})^T (\mathbf{Y} - \mathbf{U}\mathbf{S}) \right\} d\mathbf{S} \\ &\propto \text{etr} \left\{ -\frac{1}{2\sigma_n^2} \mathbf{Y}^T \mathbf{Y} + \frac{1}{2\sigma_n^2} \mathbf{Y}^T \mathbf{U} \mathbf{U}^T \mathbf{Y} \right\}. \end{aligned} \quad (12)$$

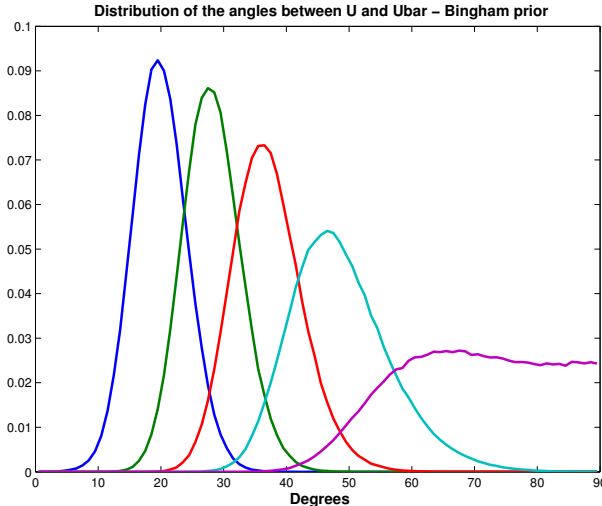


Figure 2: Distribution of the angles between  $\mathcal{R}(\mathbf{U})$  and  $\mathcal{R}(\bar{\mathbf{U}})$  for a Bingham distribution.  $N = 20$ ,  $p = 5$  and  $\kappa = 20$ .

The above distribution, along with the prior distribution  $\pi(\mathbf{U})$ , is now used to derive the MMSD estimator. We successively investigate the case of a Bingham prior and that of a vMF prior.

**Proposition 1.** *When  $\mathbf{U}$  is assigned the Bingham prior distribution, the MMSD estimator is obtained in **closed-form** as*

$$\hat{\mathbf{U}}_{mmsd-LM-B} = \mathcal{P}_p \left\{ \kappa \bar{\mathbf{U}} \bar{\mathbf{U}}^T + \frac{1}{2\sigma_n^2} \mathbf{Y} \mathbf{Y}^T \right\}. \quad (13)$$

The proof of this proposition is given in Appendix A. Therefore, in this case, the MMSD estimator has a very simple form. It consists of the principal subspace of a (weighted) combination of the a priori projection matrix  $\bar{\mathbf{U}} \bar{\mathbf{U}}^T$  and the information brought by the data through  $\mathbf{Y} \mathbf{Y}^T$ . Observe that, in this particular case of a Bingham posterior, the MMSD estimator coincides with the maximum a posteriori (MAP) estimator.

Let us now consider the case where the prior distribution of  $\mathbf{U}$  is vMF, and contrast it with the previous example. Using (12) along with along with (10), it follows that the posterior distribution now writes

$$p(\mathbf{U}|\mathbf{Y}) \propto \text{etr} \left\{ \kappa \mathbf{U}^T \bar{\mathbf{U}} + \frac{1}{2\sigma_n^2} \mathbf{U}^T \mathbf{Y} \mathbf{Y}^T \mathbf{U} \right\} \quad (14)$$

which is referred to as the Bingham-von-Mises-Fisher (BMF) distribution with parameter matrices  $\mathbf{Y} \mathbf{Y}^T$ ,  $\frac{1}{2\sigma_n^2} \mathbf{I}$  and  $\kappa \bar{\mathbf{U}}$  respectively<sup>1</sup>. Although this distribution is known [27], to our knowledge, there does not exist any analytic expression for the integral in (4) when  $\mathbf{U}|\mathbf{Y}$  has the BMF distribution (14). Therefore, the MMSD estimator cannot be computed in closed-form. In such situation, it is very common to implement a Markov chain Monte Carlo (MCMC) method for sampling according to the posterior distribution of interest. There are many MCMC algorithms that could be used for that purpose (the reader is invited to consult [32,33] for more details). However, when the full conditional distributions of the target posterior distribution can be sampled, the very popular Gibbs sampler is generally adopted for simplicity. An efficient Gibbs sampling scheme to generate random unitary matrices drawn from a BMF  $(\mathbf{A}, \mathbf{B}, \mathbf{C})$

<sup>1</sup>The matrix  $\mathbf{X}$  is said to have a BMF  $(\mathbf{A}, \mathbf{B}, \mathbf{C})$  distribution -where  $\mathbf{A}$  is an  $N \times N$  symmetric matrix,  $\mathbf{B}$  is a  $p \times p$  diagonal matrix and  $\mathbf{C}$  is an  $N \times p$  matrix- if  $p(\mathbf{X}) \propto \text{etr} \{ \mathbf{C}^T \mathbf{X} + \mathbf{B} \mathbf{X}^T \mathbf{A} \mathbf{X} \}$ .

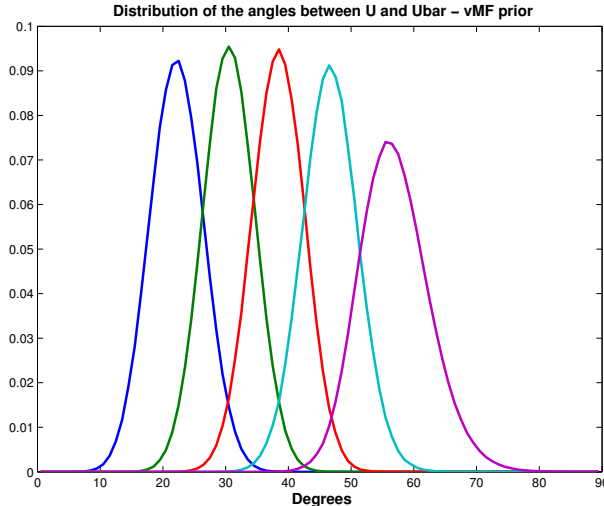


Figure 3: Distribution of the angles between  $\mathcal{R}(\mathbf{U})$  and  $\mathcal{R}(\bar{\mathbf{U}})$  for a von Mises Fisher distribution.  $N = 20$ ,  $p = 5$  and  $\kappa = 20$ .

distribution with arbitrary full-rank matrix  $\mathbf{A}$  was proposed in [34]. It amounts to sampling successively each column of  $\mathbf{U}$  by generating a random unit norm vector drawn from a (vector) BMF distribution. In our case,  $\mathbf{A} = \mathbf{Y}\mathbf{Y}^T$  whose rank is  $\min(K, N)$  and hence  $\mathbf{A}$  is rank-deficient whenever  $K < N$ , a case of most interest to us. Note also that to generate matrices  $\mathbf{U}$  drawn from the Bingham distribution in (9), we need to consider  $\mathbf{A} = \bar{\mathbf{U}}\bar{\mathbf{U}}^T$  which has rank  $p < N$ . Therefore, the scheme of [34] needs to be adapted in order to generate random matrices drawn from (14): details on how this can be achieved can be found in Appendix B. Once these matrices asymptotically distributed according to (14) are obtained, (4) can be approximated as

$$\hat{\mathbf{U}}_{\text{mmsd-LM-vMF}} \simeq \mathcal{P}_p \left\{ \frac{1}{N_r} \sum_{n=N_{\text{bi}}+1}^{N_{\text{bi}}+N_r} \mathbf{U}^{(n)} \mathbf{U}^{(n)H} \right\}. \quad (15)$$

In (15),  $N_{\text{bi}}$  is the number of burn-in samples and  $N_r$  is the number of samples used to approximate the estimator.

*Remark 3.* Interestingly enough, the above estimator in (15) is the so-called induced arithmetic mean (IAM) [35] of the set of unitary matrices  $\mathbf{U}^{(n)}$ ,  $n = N_{\text{bi}} + 1, \dots, N_{\text{bi}} + N_r$ . It differs from the Karcher mean of the set  $\mathbf{U}^{(n)}$ ,  $n = N_{\text{bi}} + 1, \dots, N_{\text{bi}} + N_r$ , which truly minimizes the sum of the distances to all  $\mathbf{U}^{(n)}$ . However, the Karcher mean may not exist and requires iterative schemes to be computed [36] while the IAM is straightforward to compute.

*Remark 4.* In the particular case where  $\mathbf{U}$  has a Bingham prior distribution, the MAP estimator of  $\mathbf{U}$  and its MMSD estimator are equal. This is no longer true when  $\mathbf{U}$  has a vMF prior distribution, and hence a BMF posterior distribution. The mode of the latter is not known in closed-form either. However, it can be approximated by selecting, among the matrices generated by the Gibbs sampler, the matrix which results in the largest value of the posterior distribution.

### 3.3 Covariance matrix model

We now consider a more complicated case where  $\mathbf{Y}$ , conditioned on  $\mathbf{U}$  and  $\mathbf{\Lambda}$ , is Gaussian distributed with zero-mean and covariance matrix

$$\mathbf{R} = \mathbb{E} \{ \mathbf{Y}\mathbf{Y}^T \} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T + \sigma_n^2 \mathbf{I} \quad (16)$$



where  $\mathbf{U}$  is an orthonormal basis for the signal subspace,  $\mathbf{\Lambda}$  is the diagonal matrix of the eigenvalues and  $\sigma_n^2$  stands for the white noise power which is assumed to be known here. As it will be more convenient and more intuitively appealing, we re-parametrize the covariance matrix as follows. The inverse of  $\mathbf{R}$  can be written as

$$\begin{aligned}\mathbf{R}^{-1} &= \mathbf{U} \left[ (\mathbf{\Lambda} + \sigma_n^2 \mathbf{I})^{-1} - \sigma_n^{-2} \mathbf{I} \right] \mathbf{U}^T + \sigma_n^{-2} \mathbf{I} \\ &= \sigma_n^{-2} \mathbf{I} - \sigma_n^{-2} \mathbf{U} \mathbf{\Lambda} (\mathbf{\Lambda} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{U}^T \\ &\triangleq \nu \mathbf{I} - \nu \mathbf{U} (\mathbf{I} - \mathbf{\Gamma}) \mathbf{U}^T\end{aligned}\tag{17}$$

where  $\nu \triangleq \sigma_n^{-2}$ ,  $\mathbf{\Gamma} \triangleq \text{diag}(\boldsymbol{\gamma})$  with  $\boldsymbol{\gamma} = [\gamma_1 \ \gamma_2 \ \cdots \ \gamma_p]^T$  and

$$0 < \gamma_k \triangleq \frac{\sigma_n^2}{\sigma_n^2 + \lambda_k} < 1.\tag{18}$$

The idea is to parametrize the problem in terms of  $\mathbf{U}$  and  $\mathbf{\Gamma}$  rather than  $\mathbf{U}$  and  $\mathbf{\Lambda}$ . The interest of this transformation is twofold. First, it enables one to express all eigenvalues with respect to the white noise level. Indeed, one has  $\mathbf{R} = \nu^{-1} \mathbf{U}_\perp \mathbf{U}_\perp^T + \nu^{-1} \mathbf{U} \mathbf{\Gamma}^{-1} \mathbf{U}^T$  where  $\mathbf{U}_\perp$  is an orthonormal basis for  $\mathcal{R}(\mathbf{U})^\perp$  and hence the  $\gamma_k$ s are representative of the scaling between the ‘‘signal’’ eigenvalues and the noise eigenvalues. In fact, they carry information about the signal-to-noise ratio since  $\gamma_k = \left(1 + \frac{\lambda_k}{\sigma_n^2}\right)^{-1}$  and  $\frac{\lambda_k}{\sigma_n^2}$  represents the SNR of the  $k$ -th signal component. Second, this new parametrization will facilitate derivation of the conditional distributions required for the Gibbs sampler.

Since  $\mathbf{Y}$  conditioned on  $\mathbf{U}$  and  $\boldsymbol{\gamma}$  is Gaussian, it follows that

$$p(\mathbf{Y}|\mathbf{U}, \boldsymbol{\gamma}) = (2\pi)^{-NK/2} |\mathbf{R}|^{-K/2} \text{etr} \left\{ -\frac{1}{2} \mathbf{Y}^T \mathbf{R}^{-1} \mathbf{Y} \right\}.\tag{19}$$

From  $\mathbf{R}^{-1} = \nu \mathbf{U}_\perp \mathbf{U}_\perp^T + \nu \mathbf{U} \mathbf{\Gamma} \mathbf{U}^T$ , it ensues that  $|\mathbf{R}^{-1}| = \nu^N |\mathbf{\Gamma}|$  and hence

$$p(\mathbf{Y}|\mathbf{U}, \boldsymbol{\gamma}) \propto |\mathbf{\Gamma}|^{K/2} \text{etr} \left\{ -\frac{1}{2} \mathbf{Y}^T [\nu \mathbf{I} - \nu \mathbf{U} (\mathbf{I} - \mathbf{\Gamma}) \mathbf{U}^T] \mathbf{Y} \right\}.\tag{20}$$

Let us now consider the prior distributions for  $\mathbf{U}$  and  $\boldsymbol{\gamma}$ . We will consider either a Bingham or vMF distribution for  $\mathbf{U}$ . As for  $\boldsymbol{\gamma}$ , we assume that  $\gamma_k$  are a priori independent random variables uniformly distributed in the interval  $[\gamma_-, \gamma_+]$ , i.e.,

$$\pi(\boldsymbol{\gamma}) = \prod_{k=1}^p (\gamma_+ - \gamma_-)^{-1} \mathbb{I}_{[\gamma_-, \gamma_+]}(\gamma_k).\tag{21}$$

The value of  $\gamma_+$  [respectively  $\gamma_-$ ] can be set to 1 [respectively 0] if a non-informative prior is desired. Otherwise, if some information is available about the SNR,  $\gamma_-$  and  $\gamma_+$  can be chosen so as to reflect this knowledge since  $\gamma_k = (1 + \text{SNR}_k)^{-1}$ :  $\gamma_+$  [resp.  $\gamma_-$ ] rules the lowest [resp. highest] value of the SNR, say  $\text{SNR}_-$  [resp.  $\text{SNR}_+$ ].

As explained in Appendix C, marginalization of  $p(\mathbf{U}, \boldsymbol{\gamma}|\mathbf{Y})$  with respect to  $\boldsymbol{\gamma}$  leads to intractable distributions  $p(\mathbf{U}|\mathbf{Y})$ . Therefore, in order to implement the MMSD estimator, we propose to draw samples from the the joint posterior distribution of  $\mathbf{U}$  and  $\boldsymbol{\gamma}$ , and then average them to estimate  $\mathbf{U}$ , similarly to what is done in (15). In Appendix C, we provide the details for generating samples according to  $p(\mathbf{U}, \boldsymbol{\gamma}|\mathbf{Y})$ , which is given by

## Bingham prior

$$\begin{aligned}
p(\mathbf{U}, \boldsymbol{\gamma} | \mathbf{Y}) &\propto p(\mathbf{Y} | \mathbf{U}, \boldsymbol{\gamma}) \pi(\mathbf{U}) \pi(\boldsymbol{\gamma}) \\
&\propto |\boldsymbol{\Gamma}|^{K/2} \left( \prod_{k=1}^p \mathbb{I}_{[\gamma_-, \gamma_+]}(\gamma_k) \right) \\
&\times \text{etr} \left\{ \kappa \mathbf{U}^T \bar{\mathbf{U}} \bar{\mathbf{U}}^T \mathbf{U} + \frac{\nu}{2} \mathbf{Y}^T \mathbf{U} (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{U}^T \mathbf{Y} \right\}.
\end{aligned} \tag{22}$$

## vMF prior

$$\begin{aligned}
p(\mathbf{U}, \boldsymbol{\gamma} | \mathbf{Y}) &\propto |\boldsymbol{\Gamma}|^{K/2} \left( \prod_{k=1}^p \mathbb{I}_{[\gamma_-, \gamma_+]}(\gamma_k) \right) \\
&\times \text{etr} \left\{ \kappa \mathbf{U}^T \bar{\mathbf{U}} + \frac{\nu}{2} \mathbf{Y}^T \mathbf{U} (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{U}^T \mathbf{Y} \right\}.
\end{aligned} \tag{23}$$

## 4 Simulations

In this section we illustrate the performance of the approach developed above through Monte Carlo simulations. In all simulations  $N = 20$ ,  $p = 5$  and  $\kappa = 20$ . The matrix  $\mathbf{S}$  is generated from a Gaussian distribution with zero-mean and covariance matrix  $\sigma_s^2 \mathbf{I}$  and the signal-to-noise ratio is defined as  $SNR = 10 \log_{10}(\sigma_s^2 / \sigma_n^2)$ . The matrix  $\mathbf{U}$  is generated from the Bingham distribution (9) or the vMF distribution (10) and, for the sake of simplicity,  $\bar{\mathbf{U}} = [\mathbf{I}_p \ \mathbf{0}]^T$ . The number of burn-in iterations in the Gibbs sampler is set to  $N_{\text{bi}} = 10$  and  $N_r = 1000$ . The MMSD estimator (4) is compared with the MAP estimator, the MMSE estimator, the usual SVD-based estimator and the estimator  $\hat{\mathbf{U}} = \bar{\mathbf{U}}$  that discards the available data and use only the a priori knowledge. The latter is referred to as ‘‘Ubar’’ in the figures. The estimators are evaluated in terms of the fraction of energy of  $\hat{\mathbf{U}}$  in  $\mathcal{R}(\mathbf{U})$ , i.e.,  $\text{AFE}(\hat{\mathbf{U}}, \mathbf{U})$ .

### 4.1 Linear model

We begin with the linear model. Figures 4 to 7 investigate the influence of  $K$  and  $SNR$  onto the performance of the estimators. Figures 4 and 5 concern the Bingham prior while the vMF prior has been used to obtain Figures 6 and 7. From inspection of these figures, the following conclusions can be drawn:

- the MMSD estimator performs better than the estimator  $\hat{\mathbf{U}} = \bar{\mathbf{U}}$ , even at low SNR. The improvement is all the more pronounced that  $K$  is large. Therefore, the MMSD estimator makes a sound use of the data to improve accuracy compared to using the prior knowledge only.
- the MMSD estimator performs better than the SVD, especially at low SNR. Moreover, and this is a distinctive feature of this Bayesian approach, it enables one to estimate the subspace even when the number of snapshots  $K$  is less than the size of the subspace  $p$ .
- for a Bingham prior, the MMSE performs very poorly since the posterior distribution of  $\mathbf{U}$  conditioned on  $\mathbf{Y}$  depends on  $\mathbf{U}\mathbf{U}^T$  only. Hence, averaging the matrix  $\mathbf{U}$  itself does not make sense, see our remark 1. In contrast, when  $\mathbf{U}$  has a vMF prior, the posterior depends on both  $\mathbf{U}$  and  $\mathbf{U}\mathbf{U}^T$ : in this case, the MMSE performs well and is close to the MMSD. Note however that the vMF prior is more restrictive than the Bingham prior.
- the MMSD estimator also outperforms the MAP estimator.

As a conclusion, the MMSD estimator performs better than most other estimators in the large majority of cases.

## 4.2 Covariance matrix model

We now conduct simulations with the covariance matrix model. The simulation parameters are essentially the same as in the previous section, except for the SNR. More precisely, the random variables  $\gamma_k$  are drawn from the uniform distribution in (21) where  $\gamma_-$  and  $\gamma_+$  are selected such that  $SNR_- = 5\text{dB}$  and  $SNR_+ = 10\text{dB}$ . The results are shown in Fig. 8 for the Bingham prior and Fig. 9 for the vMF prior. They corroborate the previous observations made on the linear model, viz that the MMSD estimator offers the best performance over all methods.

## 5 Application to hyperspectral imagery

In this section, we show how the proposed subspace estimation procedure can be efficiently used for an application to multi-band image analysis. For several decades, hyperspectral imagery has received considerable attention because of its great interest for various purposes: agriculture monitoring, mineral mapping, military concerns, etc. One of the crucial issue when analyzing such image is the spectral unmixing which aims to decompose an observed pixel  $\mathbf{y}_\ell$  into a collection of  $R = p + 1$  reference signatures,  $\mathbf{m}_1, \dots, \mathbf{m}_R$  (called *endmembers*) and to retrieve the respective proportions of these signatures (or *abundances*)  $a_{1,\ell}, \dots, a_{R,\ell}$  in this pixel [37]. To describe the physical process that links the endmembers and their abundances to the measurements, the most widely admitted mixing model is linear

$$\mathbf{y}_\ell = \sum_{r=1}^R a_{r,\ell} \mathbf{m}_r \quad (24)$$

where  $\mathbf{y}_\ell \in \mathbb{R}^N$  is the pixel spectrum measured in  $N$  spectral bands,  $\mathbf{m}_r \in \mathbb{R}^N$  ( $r = 1, \dots, R$ ) are the  $R$  endmember spectra and  $a_{r,\ell}$  ( $r = 1, \dots, R$ ) are their corresponding abundances. Due to obvious physical considerations, the abundances obey two kinds of constraints. Since they represent proportions, they must satisfy the following positivity and additivity constraints

$$\begin{cases} a_{r,\ell} \geq 0, & r = 1, \dots, R, \\ \sum_{r=1}^R a_{r,\ell} = 1. \end{cases} \quad (25)$$

Let now consider  $L$  pixels  $\mathbf{y}_1, \dots, \mathbf{y}_L$  of an hyperspectral image induced by the linear mixing model (LMM) in (24) with the abundance constraints (25). It is clear that the dataset formed by these  $L$  pixels lies in a lower-dimensional subspace  $\mathcal{U} \subset \mathbb{R}^p$ . More precisely, in this subspace  $\mathcal{U}$ , the dataset belongs to a simplex whose vertices are the endmembers  $\mathbf{m}_1, \dots, \mathbf{m}_R$  to be recovered. Most of the unmixing strategies developed in the hyperspectral imagery literature are based on this underlying geometrical formulation of the LMM. Indeed, the estimation of the endmembers is generally conducted in the lower-dimensional space  $\mathcal{U}$ , previously identified by a standard dimension reduction technique such as the principal component analysis (PCA) [37]. However, it is well known that the model linearity is a simplifying assumption and does not hold anymore in several contexts, circumventing the standard unmixing algorithms. Specifically, non-linearities are known to occur for scenes including mixtures of minerals or vegetation. As a consequence, evaluating the suitability of the LMM assumption for a given hyperspectral image is a capital question that can be conveniently addressed by the approach introduced above.

### 5.1 Synthetic data

First, we investigate the estimation of the subspace  $\mathcal{U}$  when the image pixels are non-linear functions of the abundances. For this purpose, a  $50 \times 50$  synthetic hyperspectral image is generated following a recently introduced non-linear model referred to as generalized bilinear

model (GBM). As indicated in [38], the GBM is notably well adapted to describe non-linearities due to multipath effects. It assumes that the observed pixel spectrum  $\mathbf{y}_\ell$  can be written

$$\mathbf{y}_\ell = \sum_{r=1}^R a_{r,\ell} \mathbf{m}_r + \sum_{i=1}^{R-1} \sum_{j=i+1}^R \gamma_{i,j,\ell} a_{i,\ell} a_{j,\ell} \mathbf{m}_i \odot \mathbf{m}_j \quad (26)$$

where  $\odot$  stands for the Hadamard (termwise) product and the abundances  $a_{r,\ell}$  ( $r = 1, \dots, R$ ) satisfy the constraints in (25). In (26), the parameters  $\gamma_{i,j,\ell}$  (which belong to  $[0, 1]$ ) characterize the importance of non-linear interactions between the endmembers  $\mathbf{m}_i$  and  $\mathbf{m}_j$  in the  $\ell$ -th pixel. In particular, when  $\gamma_{i,j,\ell} = 0$  ( $\forall i, j$ ), the GBM reduces to the standard LMM (24). Moreover, when  $\gamma_{i,j,\ell} = 1$  ( $\forall i, j$ ), the GBM leads to the non-linear model introduced by Fan *et al.* in [39]. In this simulation, the synthetic image has been generated using the GBM with  $R = 3$  endmember signatures extracted from a spectral library. The corresponding abundances have been uniformly drawn in the set defined by the constraints (25). We have assumed that there is no interaction between endmembers  $\mathbf{m}_1$  and  $\mathbf{m}_3$ , and between endmembers  $\mathbf{m}_2$  and  $\mathbf{m}_3$  resulting in  $\gamma_{1,3,\ell} = \gamma_{2,3,\ell} = 0$ ,  $\forall \ell$ . Moreover, the interactions between endmembers  $\mathbf{m}_1$  and  $\mathbf{m}_2$  are defined by the map of coefficients  $\gamma_{1,2,\ell}$  displayed in Fig. 10 (top, left panel) where a black (resp. white) pixel represents the lowest (resp. highest) degree of non-linearity. As can be seen in this figure, 75% of the pixels (located in the bottom and upper right squares of the image) are mixed according to the LMM resulting in  $\gamma_{1,2,\ell} = 0$ . The 25% remaining image pixels (located in the upper left square of the image) are mixed according to the GBM with nonlinearity coefficients  $\gamma_{1,2,\ell}$  radially increasing from 0 to 1 ( $\gamma_{1,2,\ell} = 0$  in the image center and  $\gamma_{1,2,\ell} = 1$  in the upper left corner of the image). Note that this image contains a majority of pixels that are mixed linearly and belong to a common subspace of  $\mathbb{R}^2$ . Conversely, the non-linearly mixed pixels do not belong to this subspace<sup>2</sup>. We propose here to estimate the local subspace  $\mathcal{U}_\ell$  where a given image pixel  $\mathbf{y}_\ell$  and its nearest spectral neighbors  $\mathcal{V}_\ell^{(K-1)}$  live ( $\mathcal{V}_\ell^{(K-1)}$  denotes the set of the  $(K-1)$ -nearest neighbors of  $\mathbf{y}_\ell$ ).

Assuming as a first approximation that all the image pixels are linearly mixed, all these pixels are approximately contained in a common 2-dimensional subspace  $\bar{\mathcal{U}}$  that can be determined by performing a PCA of  $\mathbf{y}_1, \dots, \mathbf{y}_L$  (see [40] for more details). The corresponding principal vectors spanning  $\bar{\mathcal{U}}$  are gathered in a matrix  $\bar{\mathbf{U}}$ . This matrix  $\bar{\mathbf{U}}$  is used as *a priori* knowledge regarding the 2-dimensional subspace containing  $\left\{ \mathbf{y}_\ell, \mathcal{V}_\ell^{(K-1)} \right\}_{\ell=1, \dots, L}$ . However, this crude estimation can be refined by the Bayesian estimation strategy developed in the previous sections. More precisely, for each pixel  $\mathbf{y}_\ell$ , we compute the MMSD estimator of the  $N \times p$  matrix  $\mathbf{U}_\ell$ , whose columns are supposed to span the subspace  $\mathcal{U}_\ell$  containing  $\mathbf{y}_\ell$  and its  $K-1$ -nearest neighbors  $\mathcal{V}_\ell^{(K-1)}$ . The Bayesian estimator  $\hat{\mathbf{U}}_\ell$  is computed from its closed-form expression (13), i.e., using the Bingham prior where  $\bar{\mathbf{U}}$  has been introduced above. Then, for each pixel, we evaluate the distance between the two projection matrices  $\hat{\mathbf{U}}_\ell \hat{\mathbf{U}}_\ell^T$  and  $\bar{\mathbf{U}} \bar{\mathbf{U}}^T$  onto the subspaces  $\hat{\mathcal{U}}_\ell = \mathcal{R}(\hat{\mathbf{U}}_\ell)$  and  $\bar{\mathcal{U}} = \mathcal{R}(\bar{\mathbf{U}})$ , respectively. As stated in Section 2, the natural distance between these two projection matrices is given by  $d^2(\hat{\mathbf{U}}_\ell, \bar{\mathbf{U}}) = 2(p - \text{Tr}\{\hat{\mathbf{U}}_\ell^T \bar{\mathbf{U}} \bar{\mathbf{U}}^T \hat{\mathbf{U}}_\ell\})$ . The resulting distance maps are depicted in Fig. 10 (bottom panels) for 2 non-zero values of  $\eta \triangleq 2\sigma_n^2 \kappa$  (as it can be noticed in (13), this hyperparameter  $\eta$  balances the quantity of *a priori* knowledge  $\bar{\mathbf{U}}$  included in the estimation with respect to the information brought by the data). For comparison purpose, the subspace  $\hat{\mathcal{U}}_\ell$  has been also estimated by a crude SVD of  $\left\{ \mathbf{y}_\ell, \mathcal{V}_\ell^{(K-1)} \right\}$  (top right panel). In this case,  $\hat{\mathbf{U}}_\ell$  simply reduces to the associated principal singular vectors and can be considered as the MMSD estimator of  $\mathbf{U}_\ell$  obtained for  $\eta = 0$ .

<sup>2</sup>Assuming there is a majority of image pixels that are mixed linearly is a reasonable assumption for most hyperspectral images.

These figures show that, for the 75% of the pixels generated using the LMM (bottom and right parts of the image), the subspace  $\bar{\mathcal{U}}$  estimated by an SVD of the whole dataset  $\mathbf{y}_1, \dots, \mathbf{y}_L$  is very close to the hyperplanes  $\hat{\mathcal{U}}_\ell$  locally estimated from  $\{\mathbf{y}_\ell, \mathcal{V}_\ell^{(K-1)}\}$  through the proposed approach (for any value of  $\eta$ ). Regarding the remaining 25% pixels resulting from the GBM (top left part of the image), the following comments can be made. When a crude SVD of  $\{\mathbf{y}_\ell, \mathcal{V}_\ell^{(K-1)}\}$  is conducted, i.e., when no prior knowledge is taken into account to compute the MMSD ( $\eta = 0$ , top right panel), the distance between the locally estimated subspace  $\hat{\mathcal{U}}_\ell$  and the *a priori* assumed hyperplane  $\bar{\mathcal{U}}$  does not reflect the non-linearities contained in the image. Conversely, when this crude SVD is regularized by incorporating prior knowledge with  $\eta = 0.5$  and  $\eta = 50$  (bottom left and right panels, respectively), leading to the MMSD estimator, the larger the degree of non-linearity, the larger the distance between  $\bar{\mathcal{U}}$  and  $\hat{\mathcal{U}}_\ell$ . To summarize, evaluating the distance between the MMSD estimator  $\hat{\mathcal{U}}_\ell$  and the *a priori* given matrix  $\bar{\mathcal{U}}$  allows the degree of non-linearity to be quantified. This interesting property is exploited on a real hyperspectral image in the following section.

## 5.2 Real data

The real hyperspectral image considered in this section has been acquired in 1997 over Moffett Field, CA, by the NASA spectro-imager AVIRIS. This image, depicted with composite true colors in Fig. 11 (top, left panel), has been minutely studied in [40] assuming a linear mixing model. The scene consists of a large part of a lake (black pixels, top) and a coastal area (bottom) composed of soil (brown pixels) and vegetation (green pixels), leading to  $R = 3$  endmembers whose spectra and abundance maps can be found in [40]. A simple estimation of a lower-dimensional space  $\bar{\mathcal{U}}$  where the pixels live can be conducted through a direct SVD of the whole dataset, providing the *a priori* matrix  $\bar{\mathcal{U}}$ . As in the previous section, this crude estimation can be refined by computing locally the MMSD estimators  $\hat{\mathcal{U}}_\ell$  spanning the subspaces  $\hat{\mathcal{U}}_\ell$  (bottom panels). These estimators have been also computed with  $\eta = 0$ , corresponding to an SVD of  $\{\mathbf{y}_\ell, \mathcal{V}_\ell^{(K-1)}\}$  (top, right figure). The distances between  $\bar{\mathcal{U}}$  and  $\hat{\mathcal{U}}_\ell$  have been reported in the maps of Fig. 11. Again, for  $\eta = 0$  (top, right panel), a simple local SVD is unable to locate possible non-linearities in the scene. However, for two<sup>3</sup> non-zero values  $\eta = 0.5$  and  $\eta = 50$  (bottom left and right panels, respectively), the distances between the *a priori* recovered subspace  $\bar{\mathcal{U}}$  and the MMSD-based subspace  $\hat{\mathcal{U}}_\ell$  clearly indicate that some non-linear effects occur in specific parts of the image, especially in the lake shore. Note that the non-linearities identified by the proposed algorithm are very similar to the ones highlighted in [38] where the unmixing procedure was conducted by using the GBM defined in (26). This shows the accuracy of the proposed MMSD estimator to localize the non-linearities occurring in the scene, which is interesting for the analysis of hyperspectral images.

## 6 Conclusions

This paper considered the problem of estimating a subspace using some available *a priori* information. Towards this end, a Bayesian framework was advocated, where the subspace  $\mathcal{U}$  is assumed to be drawn from an appropriate prior distribution. However, since we operate in a Grassmann manifold, the conventional MMSE approach is questionable as it amounts to minimizing a distance which is not the most meaningful on the Grassmann manifold. Consequently, we revisited the MMSE approach and proposed, as an alternative, to minimize a natural distance

---

<sup>3</sup>Additional results obtained with other values of  $\eta$  are available online at [http://dobigeon.perso.enseeiht.fr/app\\_MMSD.html](http://dobigeon.perso.enseeiht.fr/app_MMSD.html).

on the Grassmann manifold. A general framework was formulated resulting in a novel estimator which entails computing the principal eigenvectors of the posterior mean of  $\mathbf{U}\mathbf{U}^T$ . The theory was exemplified on a few simple examples, where the MMSD estimator can either be obtained in closed-form or requires resorting to an MCMC simulation method. The new approach enables one to combine efficiently the prior knowledge and the data information, resulting in a method that performs well at low SNR or with very small sample support. A successful application to the analysis of non-linearities contained in hyperspectral images was also presented.

## A Proof of Proposition 1

In this appendix, we derive the MMSD estimator for the linear model when  $\mathbf{U}$  follows a Bingham prior. In this case, the posterior distribution of  $\mathbf{U}$ , conditioned on  $\mathbf{Y}$  is given by

$$p(\mathbf{U}|\mathbf{Y}) \propto \text{etr} \left\{ \mathbf{U}^T \left[ \kappa \bar{\mathbf{U}} \bar{\mathbf{U}}^T + \frac{1}{2\sigma_n^2} \mathbf{Y} \mathbf{Y}^T \right] \mathbf{U} \right\} \quad (27)$$

which is recognized as a Bingham distribution with parameter matrix  $\kappa \bar{\mathbf{U}} \bar{\mathbf{U}}^T + \frac{1}{2\sigma_n^2} \mathbf{Y} \mathbf{Y}^T$ , i.e.,  $\mathbf{U}|\mathbf{Y} \sim \text{B} \left( \kappa \bar{\mathbf{U}} \bar{\mathbf{U}}^T + \frac{1}{2\sigma_n^2} \mathbf{Y} \mathbf{Y}^T \right)$ . Therefore, in order to derive the MMSD estimator, we need to derive the eigenvectors of  $\int \mathbf{U} \mathbf{U}^T p(\mathbf{U}|\mathbf{Y}) d\mathbf{U}$  when  $p(\mathbf{U}|\mathbf{Y})$  is a Bingham distribution. Towards this end, we make use of the following proposition.

**Proposition 2.** *Let  $\mathbf{U} \in \mathbb{R}^{N \times p}$  be an orthogonal matrix  $-\mathbf{U}^T \mathbf{U} = \mathbf{I}$ - drawn from a Bingham distribution with parameter matrix  $\mathbf{A}$*

$$p_{\text{B}}(\mathbf{U}) = \exp \{-\kappa_{\text{B}}(\mathbf{A})\} \text{etr} \{ \mathbf{U}^T \mathbf{A} \mathbf{U} \} \quad (28)$$

with  $\kappa_{\text{B}}(\mathbf{A}) = \ln {}_1F_1 \left( \frac{1}{2}p, \frac{1}{2}N; \mathbf{A} \right)$ . Let  $\mathbf{A} = \mathbf{U}_a \boldsymbol{\Lambda}_a \mathbf{U}_a^T$  denote the eigenvalue decomposition of  $\mathbf{A}$  where the eigenvalues are ordered in descending order. Let us define  $\mathbf{M} = \int \mathbf{U} \mathbf{U}^T p_{\text{B}}(\mathbf{U}) d\mathbf{U}$ . Then the eigenvalue decomposition of  $\mathbf{M}$  writes

$$\mathbf{M} = \exp \{-\kappa_{\text{B}}(\mathbf{A})\} \mathbf{U}_a \boldsymbol{\Gamma} \mathbf{U}_a^T$$

with  $\boldsymbol{\Gamma} = \frac{\partial \exp\{\kappa_{\text{B}}(\mathbf{A})\}}{\partial \boldsymbol{\Lambda}_a}$  and  $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_N$  where  $\gamma_n = \boldsymbol{\Gamma}(n, n)$ .

*Proof.* For notational convenience, let us work with the projection matrix  $\mathbf{P} = \mathbf{U} \mathbf{U}^T$  whose distribution on the Grassmann manifold is [27]

$$p(\mathbf{P}) = \exp \{-\kappa_{\text{B}}(\mathbf{A})\} \text{etr} \{ \mathbf{P} \mathbf{A} \}. \quad (29)$$

We have then that

$$\begin{aligned} \mathbf{M} &= \exp \{-\kappa_{\text{B}}(\mathbf{A})\} \int \mathbf{P} \text{etr} \{ \mathbf{P} \mathbf{U}_a \boldsymbol{\Lambda}_a \mathbf{U}_a^T \} d\mathbf{P} \\ &= \exp \{-\kappa_{\text{B}}(\mathbf{A})\} \mathbf{U}_a \left[ \int \mathbf{U}_a^T \mathbf{P} \mathbf{U}_a \text{etr} \{ \mathbf{U}_a^T \mathbf{P} \mathbf{U}_a \boldsymbol{\Lambda}_a \} d\mathbf{P} \right] \mathbf{U}_a^T \\ &= \exp \{-\kappa_{\text{B}}(\mathbf{A})\} \mathbf{U}_a \left[ \int \mathbf{P} \text{etr} \{ \mathbf{P} \boldsymbol{\Lambda}_a \} d\mathbf{P} \right] \mathbf{U}_a^T \\ &= \exp \{-\kappa_{\text{B}}(\mathbf{A})\} \mathbf{U}_a \boldsymbol{\Gamma} \mathbf{U}_a^T. \end{aligned}$$

Moreover  $\boldsymbol{\Gamma}$  is diagonal since, for any orthogonal diagonal matrix  $\mathbf{D}$ ,

$$\begin{aligned} \boldsymbol{\Gamma} \mathbf{D} &= \int \mathbf{P} \mathbf{D} \text{etr} \{ \mathbf{P} \boldsymbol{\Lambda}_a \} d\mathbf{P} \\ &= \mathbf{D} \left[ \int \mathbf{D}^T \mathbf{P} \mathbf{D} \text{etr} \{ \mathbf{D}^T \mathbf{P} \mathbf{D} \boldsymbol{\Lambda}_a \mathbf{D} \} d\mathbf{P} \right] \\ &= \mathbf{D} \int \mathbf{P} \text{etr} \{ \mathbf{P} \boldsymbol{\Lambda}_a \} d\mathbf{P} \\ &= \mathbf{D} \boldsymbol{\Gamma} \end{aligned}$$

where, to obtain the third line, we made use of the fact that  $\mathbf{D}^T \boldsymbol{\Lambda}_a \mathbf{D} = \boldsymbol{\Lambda}_a$ . It follows that the eigenvectors of  $\mathbf{M}$  and  $\mathbf{A}$  coincide, and that the eigenvalues of  $\mathbf{M}$  are  $\exp \{-\kappa_{\text{B}}(\mathbf{A})\} \gamma_n$ , for

$n = 1, \dots, N$ . Moreover, it is known that  $\exp\{-\kappa_{\mathbf{B}}(\mathbf{A})\} = \exp\{-\kappa_{\mathbf{B}}(\mathbf{\Lambda}_a)\}$  and, from (29), one has

$$\exp\{\kappa_{\mathbf{B}}(\mathbf{\Lambda}_a)\} = \int \text{etr}\{\mathbf{P}\mathbf{\Lambda}_a\} d\mathbf{P}.$$

Differentiating the latter equation with respect to  $\lambda_a(k)$  and denoting  $p_n = \mathbf{P}(n, n)$ , one obtains

$$\begin{aligned} \frac{\partial \exp\{\kappa_{\mathbf{B}}(\mathbf{\Lambda}_a)\}}{\partial \lambda_a(k)} &= \frac{\partial}{\partial \lambda_a(k)} \int \exp\left\{\sum_{n=1}^N \lambda_a(n) p_n\right\} d\mathbf{P} \\ &= \int p_k \text{etr}\{\mathbf{P}\mathbf{\Lambda}_a\} d\mathbf{P} \\ &= \gamma_k. \end{aligned}$$

The previous equation enables one to relate the eigenvalues of  $\mathbf{A}$  and those of  $\mathbf{M}$ . It remains to prove that  $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_N$ . Towards this end, we make use of a very general theorem due to Letac [41], which is briefly outlined below. Let

$$P(\mu, \mathbf{A})(d\mathbf{X}) = \exp\{\kappa_{\mu}(\mathbf{A})\} \text{etr}\{\mathbf{X}^T \mathbf{A}\} \mu(d\mathbf{X})$$

be a probability associated with a unitarily invariant measure  $\mu$  on the set of  $N \times N$  symmetric matrices. Consider the case of a diagonal matrix  $\mathbf{A} = \text{diag}(a_1, a_2, \dots, a_N)$  with  $a_1 \geq a_2 \geq \dots \geq a_N$ . Then [41] proves that  $\mathbf{M} = \int \mathbf{X} P(\mu, \mathbf{A})(d\mathbf{X})$  is also diagonal, and moreover if  $\mathbf{M} = \text{diag}(m_1, m_2, \dots, m_N)$  then  $m_1 \geq m_2 \geq \dots \geq m_N$ . Use of this theorem completes the proof of the proposition.  $\square$

It then follows that the posterior distribution is  $\mathbf{U}|\mathbf{Y} \sim \text{B}\left(\kappa \bar{\mathbf{U}}\bar{\mathbf{U}}^T + \frac{1}{2\sigma_n^2} \mathbf{Y}\mathbf{Y}^T\right)$  when  $\mathbf{U}$  has a Bingham prior. Thus, the eigenvectors of  $\int \mathbf{U}\mathbf{U}^T p(\mathbf{U}|\mathbf{Y}) d\mathbf{U}$  coincide with those of  $\kappa \bar{\mathbf{U}}\bar{\mathbf{U}}^T + \frac{1}{2\sigma_n^2} \mathbf{Y}\mathbf{Y}^T$ , with the same ordering of their eigenvalues. Consequently, the MMSD estimator is given by (13), which concludes the proof.

## B Sampling from the Bingham-von Mises Fisher distribution

In this appendix, we show how to sample a unitary random matrix  $\mathbf{X} \in \mathbb{R}^{N \times p}$  from a (matrix) Bingham von Mises Fisher (BMF) distribution,  $\mathbf{X} \sim \text{BMF}(\mathbf{A}, \mathbf{B}, \mathbf{C})$ . As will be explained shortly, this amounts to sampling successively each column of  $\mathbf{X}$ , and entails generating a random unit norm vector drawn from a (vector) BMF distribution. We briefly review how to sample the columns of  $\mathbf{X}$  and then explain how to sample from a vector BMF distribution.

### B.1 The matrix BMF distribution

The density of  $\mathbf{X} \sim \text{BMF}(\mathbf{A}, \mathbf{B}, \mathbf{C})$  is given by

$$\begin{aligned} p(\mathbf{X}|\mathbf{A}, \mathbf{B}, \mathbf{C}) &\propto \text{etr}\{\mathbf{C}^T \mathbf{X} + \mathbf{B}\mathbf{X}^T \mathbf{A}\mathbf{X}\} \\ &\propto \prod_{k=1}^p \exp\{\mathbf{c}_k^T \mathbf{x}_k + \mathbf{B}(k, k) \mathbf{x}_k^T \mathbf{A} \mathbf{x}_k\} \end{aligned} \quad (30)$$

where  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_p]$  and  $\mathbf{C} = [\mathbf{c}_1 \ \mathbf{c}_2 \ \dots \ \mathbf{c}_p]$ . In [34] a Gibbs-sampling strategy was presented in order to sample from this distribution, in the case where  $\mathbf{A}$  is full-rank. We consider here a situation where  $\mathbf{A}$  is rank-deficient and therefore we need to bring appropriate modifications to the scheme of [34] in order to handle the rank deficiency of  $\mathbf{A}$ . As evidenced from (30)



the distribution of  $\mathbf{A}$  is a product of vector BMF distributions, except that the columns of  $\mathbf{X}$  are not statistically independent since they are orthogonal with probability one. Let us rewrite  $\mathbf{X}$  as  $\mathbf{X} = [\mathbf{x}_1 \ \cdots \ \mathbf{x}_{k-1} \ \mathbf{Q}_\perp \mathbf{z} \ \mathbf{x}_{k+1} \ \cdots \ \mathbf{x}_p]$  where  $\mathbf{z} \in \mathcal{S}_{N-p+1} = \{\mathbf{x} \in \mathbb{R}^{N-p+1 \times 1}; \mathbf{x}^T \mathbf{x} = 1\}$  and  $\mathbf{Q}_\perp$  is an  $N \times N - p + 1$  orthonormal basis for  $\mathcal{R}(\mathbf{X}_{-k})^\perp$  where  $\mathbf{X}_{-k}$  stands for the matrix  $\mathbf{X}$  with its  $k$ -th column removed. As shown in [34] the conditional density of  $\mathbf{z}$  given  $\mathbf{X}_{-k}$  is

$$\begin{aligned} p(\mathbf{z}|\mathbf{X}_{-k}) &\propto \exp\{\mathbf{c}_k^T \mathbf{Q}_\perp \mathbf{z} + \mathbf{B}(k, k) \mathbf{z}^T \mathbf{Q}_\perp^T \mathbf{A} \mathbf{Q}_\perp \mathbf{z}\} \\ &\propto \exp\{\tilde{\mathbf{c}}_k^T \mathbf{z} + \mathbf{z}^T \tilde{\mathbf{A}} \mathbf{z}\} \end{aligned} \quad (31)$$

where  $\tilde{\mathbf{c}}_k = \mathbf{Q}_\perp^T \mathbf{c}_k$  and  $\tilde{\mathbf{A}} = \mathbf{B}(k, k) \mathbf{Q}_\perp^T \mathbf{A} \mathbf{Q}_\perp$ . Therefore,  $\mathbf{z}|\mathbf{X}_{-k}$  follows a vector BMF distribution  $\mathbf{z}|\mathbf{X}_{-k} \sim \text{vBMF}(\tilde{\mathbf{A}}, \tilde{\mathbf{c}}_k)$ . A Markov chain that converges to BMF( $\mathbf{A}, \mathbf{B}, \mathbf{C}$ ) can thus be constructed as follows:

**Input:** initial value  $\mathbf{X}^{(0)}$

- 1: **for**  $k = 1, \dots, p$  (random order) **do**
- 2:   compute a basis  $\mathbf{Q}_\perp$  for the null space of  $\mathbf{X}_{-k}$  and set  $\mathbf{z} = \mathbf{Q}_\perp^T \mathbf{x}_k$ .
- 3:   compute  $\tilde{\mathbf{c}}_k = \mathbf{Q}_\perp^T \mathbf{c}_k$  and  $\tilde{\mathbf{A}} = \mathbf{B}(k, k) \mathbf{Q}_\perp^T \mathbf{A} \mathbf{Q}_\perp$ .
- 4:   sample  $\mathbf{z}$  from a vBMF( $\tilde{\mathbf{A}}, \tilde{\mathbf{c}}_k$ ) distribution (*see next section*).
- 5:   set  $\mathbf{x}_k = \mathbf{Q}_\perp \mathbf{z}$ .
- 6: **end for**

## B.2 The vector BMF distribution

The core part of the above algorithm, see line 4, is to draw a unit-norm random vector  $\mathbf{x}$  distributed according to a vector Bingham-von Mises Fisher distribution. The latter distribution on the  $M$ - dimensional sphere has a density with respect to the uniform distribution given by

$$p(\mathbf{x}|\mathbf{c}, \mathbf{A}) \propto \exp\{\mathbf{c}^T \mathbf{x} + \mathbf{x}^T \mathbf{A} \mathbf{x}\}, \mathbf{x} \in \mathcal{S}_M. \quad (32)$$

In [34] a Gibbs-sampling strategy was presented in order to sample from this distribution. While  $\mathbf{A}$  was assumed to be full-rank in [34], we consider here a situation where  $\mathbf{A}$  is rank-deficient, i.e. its eigenvalue decomposition can be written as  $\mathbf{A} = \mathbf{E} \mathbf{\Lambda} \mathbf{E}^T$  where  $\mathbf{E}$  stands for the orthonormal matrix of the eigenvectors and  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_r, 0, \dots, 0)$  is the diagonal matrix of its eigenvalues. Our derivation follows along the same lines as in [34] with the appropriate modifications due to the rank deficiency of  $\mathbf{A}$ . Let  $\mathbf{y} = \mathbf{E}^T \mathbf{x} \in \mathcal{S}_M$  and  $\mathbf{d} = \mathbf{E}^T \mathbf{c}$ . Since  $y_M^2 = 1 - \sum_{k=1}^{M-1} y_k^2$ , the uniform density in terms of the unconstrained coordinates  $\{y_1, y_2, \dots, y_{M-1}\}$  is proportional to  $|y_M|^{-1}$  and the density of  $\{y_1, y_2, \dots, y_{M-1}\}$  is given by [34]

$$\begin{aligned} p(\mathbf{y}|\mathbf{d}, \mathbf{E}) &\propto \exp\{\mathbf{d}^T \mathbf{y} + \mathbf{y}^T \mathbf{\Lambda} \mathbf{y}\} |y_M|^{-1}, \quad y_M^2 = 1 - \sum_{k=1}^{M-1} y_k^2 \\ &\propto \exp\left\{\sum_{k=1}^M d_k y_k + \sum_{k=1}^r \lambda_k y_k^2\right\} |y_M|^{-1}. \end{aligned} \quad (33)$$

In order to sample from this distribution, a Gibbs sampling strategy is advocated. Towards this end, we need to derive the conditional distributions of  $y_k$ , given  $\mathbf{y}_{-k}$  where  $\mathbf{y}_{-k}$  stands for the vector  $\mathbf{y}$  with its  $k$ -th component removed. Similarly to [34], let us make the change of variables  $\theta_k = y_k^2$  and let  $\mathbf{q} = \left[ \frac{y_1^2}{1-y_k^2} \quad \frac{y_2^2}{1-y_k^2} \quad \cdots \quad \frac{y_M^2}{1-y_k^2} \right]^T$ , so that  $\{y_1^2, y_2^2, \dots, y_M^2\} = \{\theta_k, (1 - \theta_k) \mathbf{q}_{-k}\}$ . Since this change of variables is not bijective, i.e.  $y_k \pm \theta_k^{1/2}$ , we need to introduce the sign  $s_k$  of  $y_k$ , and we let  $\mathbf{s} = [s_1 \ s_2 \ \cdots \ s_M]^T$ . Note that  $y_M^2 = 1 - \sum_{k=1}^{M-1} y_k^2$ ,

$|y_M| = (1 - \theta_k)^{1/2} q_M^{1/2}$  and  $q_M = 1 - \sum_{\ell=1, \ell \neq k}^{M-1} q_\ell$ . As shown in [34], the Jacobian of the transformation from  $\{y_1, y_2, \dots, y_{M-1}\}$  to  $\{\theta, q_1, \dots, q_{k-1}, q_{k+1}, \dots, q_{M-1}\}$  is proportional to  $\theta_k^{-1/2} (1 - \theta_k)^{(M-2)/2} \prod_{\ell=1, \ell \neq k}^{M-1} q_\ell^{-1/2}$ , and therefore the joint distribution of  $\theta_k, s_k, \mathbf{q}_{-k}, \mathbf{s}_{-k}$  can be written as

$$\begin{aligned}
p(\theta_k, s_k, \mathbf{q}_{-k}, \mathbf{s}_{-k}) &\propto \theta_k^{-1/2} (1 - \theta_k)^{(M-3)/2} \left( \prod_{\ell \neq k} q_\ell^{-1/2} \right) \\
&\times \exp \left\{ s_k \theta_k^{1/2} d_k + (1 - \theta_k)^{1/2} \sum_{\ell \neq k} d_\ell s_\ell q_\ell^{1/2} \right\} \\
&\times \begin{cases} \exp \left\{ \theta_k \lambda_k + (1 - \theta_k) \sum_{\ell=1, \ell \neq k}^r q_\ell \lambda_\ell \right\} & 1 \leq k \leq r \\ \exp \left\{ (1 - \theta_k) \sum_{\ell=1}^r q_\ell \lambda_\ell \right\} & r+1 \leq k \leq M \end{cases}. \quad (34)
\end{aligned}$$

It follows that

- for  $k \in [1, r]$

$$\begin{aligned}
p(\theta_k, s_k | \mathbf{q}_{-k}, \mathbf{s}_{-k}) &\propto \theta_k^{-1/2} (1 - \theta_k)^{(M-3)/2} \exp \left\{ \theta_k \lambda_k + (1 - \theta_k) \mathbf{q}_{-k}^T \boldsymbol{\lambda}_{-k} \right\} \\
&\times \exp \left\{ s_k \theta_k^{1/2} d_k + (1 - \theta_k)^{1/2} \left[ \mathbf{s}_{-k} \odot \mathbf{q}_{-k}^{1/2} \right]^T \mathbf{d}_{-k} \right\}. \quad (35)
\end{aligned}$$

- for  $k \in [r+1, M]$

$$\begin{aligned}
p(\theta_k, s_k | \mathbf{q}_{-k}, \mathbf{s}_{-k}) &\propto \theta_k^{-1/2} (1 - \theta_k)^{(M-3)/2} \exp \left\{ (1 - \theta_k) \mathbf{q}^T \boldsymbol{\lambda} \right\} \\
&\times \exp \left\{ s_k \theta_k^{1/2} d_k + (1 - \theta_k)^{1/2} \left[ \mathbf{s}_{-k} \odot \mathbf{q}_{-k}^{1/2} \right]^T \mathbf{d}_{-k} \right\}. \quad (36)
\end{aligned}$$

In the previous equations,  $\odot$  stands for the element-wise vector or matrix product and  $\mathbf{q}_{-k}^{1/2}$  is a short-hand notation to designate the vector  $\left[ q_1^{1/2} \dots q_{k-1}^{1/2} q_{k+1}^{1/2} \dots q_M^{1/2} \right]^T$ . In order to sample from  $p(\theta_k, s_k | \mathbf{q}_{-k}, \mathbf{s}_{-k})$ , we first sample  $\theta_k$  from

$$\begin{aligned}
p(\theta_k | \mathbf{q}_{-k}, \mathbf{s}_{-k}) &= p(\theta_k, s_k = -1 | \mathbf{q}_{-k}, \mathbf{s}_{-k}) + p(\theta_k, s_k = 1 | \mathbf{q}_{-k}, \mathbf{s}_{-k}) \\
&\propto \theta_k^{-1/2} (1 - \theta_k)^{(M-3)/2} \exp \left\{ a_k \theta_k + b_k (1 - \theta_k)^{1/2} \right\} \\
&\times \left[ \exp \left\{ -d_k \theta_k^{1/2} \right\} + \exp \left\{ -d_k \theta_k^{1/2} \right\} \right] \quad (37)
\end{aligned}$$

where  $b_k = \left[ \mathbf{s}_{-k} \odot \mathbf{q}_{-k}^{1/2} \right]^T \mathbf{d}_{-k}$  and

$$a_k = \begin{cases} \lambda_k - \mathbf{q}_{-k}^T \boldsymbol{\lambda}_{-k} & k \in [1, r] \\ -\mathbf{q}^T \boldsymbol{\lambda} & k \in [r+1, M] \end{cases}. \quad (38)$$

Next, we sample  $s_k \in \{-1, +1\}$  with probabilities proportional to  $\left( e^{-d_k \theta_k^{1/2}}, e^{+d_k \theta_k^{1/2}} \right)$ . In order to sample from the distribution in (37), an efficient rejection sampling scheme was proposed in [34], where the proposal distribution is a beta distribution with suitably chosen parameters.

## C MCMC implementation of the MMSD estimator in the covariance matrix model

In this appendix, we provide the necessary details for MCMC implementation of the MMSD estimator in the covariance matrix model. We successively investigate the case of a Bingham prior and the case of a vMF prior.

### C.1 Bingham prior

When the prior distribution of  $\mathbf{U}$  is the Bingham distribution of (9), the joint posterior distribution of  $\mathbf{U}$  and  $\boldsymbol{\gamma}$  is given by (22). Since the MMSD estimator involves the posterior distribution  $p(\mathbf{U}|\mathbf{Y})$  of  $\mathbf{U}$  only, a natural way to proceed is to marginalize (22) with respect to  $\boldsymbol{\gamma}$ . Let  $\mathbf{Z} = \mathbf{Y}^T \mathbf{U} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \cdots \ \mathbf{z}_p]$ . Then, from (22) one has

$$\begin{aligned}
p(\mathbf{U}|\mathbf{Y}) &= \int p(\mathbf{U}, \boldsymbol{\gamma}|\mathbf{Y}) d\boldsymbol{\gamma} \\
&\propto \text{etr} \left\{ \kappa \mathbf{U}^T \bar{\mathbf{U}} \bar{\mathbf{U}}^T \mathbf{U} + \frac{\nu}{2} \mathbf{U}^T \mathbf{Y} \mathbf{Y}^T \mathbf{U} \right\} \\
&\times \prod_{k=1}^p \int_{\gamma_-}^{\gamma_+} \gamma_k^{K/2} \exp \left\{ -\frac{\nu}{2} \gamma_k \|\mathbf{z}_k\|^2 \right\} d\gamma_k \\
&\propto \text{etr} \left\{ \kappa \mathbf{U}^T \bar{\mathbf{U}} \bar{\mathbf{U}}^T \mathbf{U} + \frac{\nu}{2} \mathbf{U}^T \mathbf{Y} \mathbf{Y}^T \mathbf{U} \right\} \\
&\times \prod_{k=1}^p \|\mathbf{z}_k\|^{-2(1+K/2)} \left[ \gamma \left( \frac{\nu}{2} \gamma_+ \|\mathbf{z}_k\|^2, 1 + \frac{K}{2} \right) - \gamma \left( \frac{\nu}{2} \gamma_- \|\mathbf{z}_k\|^2, 1 + \frac{K}{2} \right) \right] \quad (39)
\end{aligned}$$

where  $\gamma(x, a) = \int_0^x t^{a-1} e^{-t} dt$  is the incomplete Gamma function. Unfortunately, the above distribution does not belong to any known family and it is thus problematic to generate samples drawn from it. Instead, in order to sample according to (22), we propose to use a Gibbs sampler drawing samples according to  $p(\mathbf{U}|\mathbf{Y}, \boldsymbol{\gamma})$  and  $p(\gamma_k|\mathbf{Y}, \mathbf{U})$  for  $k = 1, \dots, p$ . From (22), the conditional distribution of  $\mathbf{U}$  is

$$p(\mathbf{U}|\mathbf{Y}, \boldsymbol{\gamma}) \propto \text{etr} \left\{ \kappa \mathbf{U}^T \bar{\mathbf{U}} \bar{\mathbf{U}}^T \mathbf{U} + \frac{\nu}{2} (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{U}^T \mathbf{Y} \mathbf{Y}^T \mathbf{U} \right\} \quad (40)$$

which is recognized as a (modified) Bingham distribution<sup>4</sup>

$$\mathbf{U}|\mathbf{Y}, \boldsymbol{\gamma} \sim \tilde{\mathbf{B}} \left( \bar{\mathbf{U}} \bar{\mathbf{U}}^T, \kappa \mathbf{I}, \mathbf{Y} \mathbf{Y}^T, \frac{\nu}{2} (\mathbf{I} - \boldsymbol{\Gamma}) \right). \quad (41)$$

Let us now turn to the conditional distribution of  $\boldsymbol{\gamma}|\mathbf{Y}, \mathbf{U}$ . From (22) one has

$$\begin{aligned}
p(\boldsymbol{\gamma}|\mathbf{Y}, \mathbf{U}) &\propto |\boldsymbol{\Gamma}|^{K/2} \text{etr} \left\{ -\frac{\nu}{2} \mathbf{Z} \boldsymbol{\Gamma} \mathbf{Z}^T \right\} \left( \prod_{k=1}^p \mathbb{I}_{[\gamma_-, \gamma_+]}(\gamma_k) \right) \\
&\propto \prod_{k=1}^p \left[ \gamma_k^{K/2} \exp \left\{ -\frac{\nu}{2} \|\mathbf{z}_k\|^2 \gamma_k \right\} \mathbb{I}_{[\gamma_-, \gamma_+]}(\gamma_k) \right] \quad (42)
\end{aligned}$$

which is the product of independent gamma distributions with parameters  $\frac{K}{2} + 1$  and  $\frac{\nu}{2} \|\mathbf{z}_k\|^2$ , truncated in the interval  $[\gamma_-, \gamma_+]$ . We denote this distribution as  $\gamma_k \sim \mathcal{G}_t \left( \frac{K}{2} + 1, \frac{\nu}{2} \|\mathbf{z}_k\|^2, \gamma_-, \gamma_+ \right)$ . Random variables with such a distribution can be efficiently generated using the accept-reject scheme of [42].

---

<sup>4</sup> $\mathbf{X} \sim \tilde{\mathbf{B}}(\mathbf{A}_1, \mathbf{B}_1, \mathbf{A}_2, \mathbf{B}_2) \Leftrightarrow p(\mathbf{X}) \propto \text{etr} \{ \mathbf{B}_1 \mathbf{X}^T \mathbf{A}_1 \mathbf{X} + \mathbf{B}_2 \mathbf{X}^T \mathbf{A}_2 \mathbf{X} \}$

The above conditional distributions can now be used in a Gibbs sampler, as described in Algorithm 1. The so-generated matrices  $\mathbf{U}^{(n)}$  can be used similarly to (15) to obtain the MMSD estimator.

---

**Algorithm 1** Covariance matrix model. Gibbs sampler for the Bingham prior distribution.

---

**Input:** initial values  $\mathbf{U}^{(0)}, \boldsymbol{\gamma}^{(0)}$

1: **for**  $n = 1, \dots, N_{bi} + N_r$  **do**

2: sample  $\mathbf{U}^{(n)}$  from  $\tilde{\mathbf{B}}\left(\kappa\mathbf{I}, \bar{\mathbf{U}}\bar{\mathbf{U}}^T, \frac{\nu}{2}\left(\mathbf{I} - \boldsymbol{\Gamma}^{(n-1)}\right), \mathbf{Y}\mathbf{Y}^T\right)$  in (40).

3: for  $k = 1, \dots, p$ , sample  $\gamma_k^{(n)}$  from  $\mathcal{G}_t\left(\frac{K}{2} + 1, \frac{\nu}{2}\left\|\mathbf{Y}^T \mathbf{u}_k^{(n)}\right\|^2, \gamma_-, \gamma_+\right)$  in (42).

4: **end for**

**Output:** sequence of random variables  $\mathbf{U}^{(n)}$  and  $\boldsymbol{\gamma}^{(n)}$

---

## C.2 von Mises Fisher prior

When  $\mathbf{U}$  has a vMF prior distribution, the joint posterior distribution of  $\mathbf{U}$  and  $\boldsymbol{\gamma}$  is given by (23). Marginalizing the latter with respect to  $\boldsymbol{\gamma}$  will yield a posterior distribution  $p(\mathbf{U}|\mathbf{Y})$  similar to that of equation (39), except that the term  $\kappa\mathbf{U}^T\bar{\mathbf{U}}\bar{\mathbf{U}}^T\mathbf{U}$  should be replaced by  $\kappa\mathbf{U}^T\bar{\mathbf{U}}$ . Again this leads to an intractable posterior. Therefore, as done previously, we consider drawing samples according to  $p(\mathbf{U}|\mathbf{Y}, \boldsymbol{\gamma})$  and  $p(\boldsymbol{\gamma}|\mathbf{Y}, \mathbf{U})$ . The latter distribution will still be given by (42) while  $p(\mathbf{U}|\mathbf{Y}, \boldsymbol{\gamma})$  now takes the form

$$p(\mathbf{U}|\mathbf{Y}, \boldsymbol{\gamma}) \propto \text{etr}\left\{\kappa\mathbf{U}^T\bar{\mathbf{U}} + \frac{\nu}{2}\left(\mathbf{I} - \boldsymbol{\Gamma}\right)\mathbf{U}^T\mathbf{Y}\mathbf{Y}^T\mathbf{U}\right\} \quad (43)$$

which is recognized as a BMF distribution  $\mathbf{U}|\mathbf{Y}, \boldsymbol{\gamma} \sim \text{BMF}\left(\mathbf{Y}\mathbf{Y}^T, \frac{\nu}{2}\left(\mathbf{I} - \boldsymbol{\Gamma}\right), \kappa\bar{\mathbf{U}}\right)$ . Therefore only line 2 of the Gibbs sampler in Table 1 needs to be modified, which yields the Gibbs sampler of Algorithm 2 in the case of a vMF prior.

---

**Algorithm 2** Covariance matrix model. Gibbs sampler for the vMF prior distribution.

---

**Input:** initial values  $\mathbf{U}^{(0)}, \boldsymbol{\gamma}^{(0)}$

1: **for**  $n = 1, \dots, N_{bi} + N_r$  **do**

2: sample  $\mathbf{U}^{(n)}$  from  $\text{BMF}\left(\mathbf{Y}\mathbf{Y}^T, \frac{\nu}{2}\left(\mathbf{I} - \boldsymbol{\Gamma}^{(n-1)}\right), \kappa\bar{\mathbf{U}}\right)$  in (43).

3: for  $k = 1, \dots, p$ , sample  $\gamma_k^{(n)}$  from  $\mathcal{G}_t\left(\frac{K}{2} + 1, \frac{\nu}{2}\left\|\mathbf{Y}^T \mathbf{u}_k^{(n)}\right\|^2, \gamma_-, \gamma_+\right)$  in (42).

4: **end for**

**Output:** sequence of random variables  $\mathbf{U}^{(n)}$  and  $\boldsymbol{\gamma}^{(n)}$

---

## Acknowledgment

The authors would like to thank Prof. Kit Bingham from the University of Minnesota for insightful comments on the Bingham distribution and for pointing reference [43]. They are also indebted to Prof. Gérard Letac, University of Toulouse, for fruitful discussions leading to the proof of Proposition 2 given in Appendix A.

## References

- [1] L. L. Scharf, *Statistical Signal Processing: Detection, Estimation and Time Series Analysis*. Reading, MA: Addison Wesley, 1991.
- [2] R. O. Schmidt, “Multiple emitter location and signal parameter estimation,” *IEEE Transactions Antennas Propagation*, vol. 34, no. 3, pp. 276–280, March 1986.
- [3] R. Roy and T. Kailath, “ESPRIT- estimation of signal parameters via rotational invariance techniques,” *IEEE Transactions Acoustics Speech Signal Processing*, vol. 37, no. 7, pp. 984–995, July 1989.
- [4] R. Kumaresan and D. Tufts, “Estimating the parameters of exponentially damped sinusoids and pole-zero modeling in noise,” *IEEE Transactions Acoustics Speech Signal Processing*, vol. 30, no. 6, pp. 833–840, December 1982.
- [5] —, “Estimating the angles of arrival of multiple plane waves,” *IEEE Transactions Aerospace Electronic Systems*, vol. 19, no. 1, pp. 134–139, January 1983.
- [6] B. Ottersten, M. Viberg, P. Stoica, and A. Nehorai, “Exact and large sample maximum likelihood techniques for parameter estimation and detection in array processing,” in *Radar Array Processing*, S. Haykin, J. Litva, and T. Shepherd, Eds. Berlin: Springer Verlag, 1993, ch. 4, pp. 99–151.
- [7] P. Stoica and A. Nehorai, “MUSIC, maximum likelihood and Cramér-Rao bound,” *IEEE Transactions Acoustics Speech Signal Processing*, vol. 37, no. 5, pp. 720–741, May 1989.
- [8] —, “MUSIC, maximum likelihood and Cramér-Rao bound: Further results and comparisons,” *IEEE Transactions Acoustics Speech Signal Processing*, vol. 38, no. 12, pp. 2140–2150, December 1990.
- [9] J. Thomas, L. Scharf, and D. Tufts, “The probability of a subspace swap in the SVD,” *IEEE Transactions Signal Processing*, vol. 43, no. 3, pp. 730–736, March 1995.
- [10] M. Hawkes, A. Nehorai, and P. Stoica, “Performance breakdown of subspace-based methods: prediction and cure,” in *Proceedings ICASSP*, May 2001, pp. 4005–4008.
- [11] B. Johnson, Y. Abramovich, and X. Mestre, “The role of subspace swap in MUSIC performance breakdown,” in *Proceedings ICASSP*, March 2008, pp. 2473–2476.
- [12] R. R. Nadakuditi and F. Benaych-Georges, “The breakdown point of signal subspace estimation,” in *Proceedings SAM*, Israel, 4-7 October 2010, pp. 177–180.
- [13] X. Mestre, “Improved estimation of eigenvalues and eigenvectors of covariance matrices using their sample estimates,” *IEEE Transactions Information Theory*, vol. 54, no. 11, pp. 5113–5129, November 2008.
- [14] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Englewood Cliffs, NJ: Prentice Hall, 1993.
- [15] A. Edelman, T. Arias, and S. Smith, “The geometry of algorithms with orthogonality constraints,” *SIAM Journal Matrix Analysis Applications*, vol. 20, no. 2, pp. 303–353, 1998.
- [16] G. Golub and C. V. Loan, *Matrix Computations*, 3rd ed. Baltimore: John Hopkins University Press, 1996.

- [17] A. Srivastava, “A Bayesian approach to geometric subspace estimation,” *IEEE Transactions Signal Processing*, vol. 48, no. 5, pp. 1390–1400, May 2000.
- [18] U. Grenander, M. I. Miller, and A. Srivastava, “Hilbert-Schmidt lower bounds for estimators on matrix Lie groups for ATR,” *IEEE Transactions Pattern Analysis Machine Intelligence*, vol. 20, no. 8, pp. 790–802, August 1998.
- [19] Q. Rentmeesters, P. Absil, P. Van Dooren, K. Gallivan, and A. Srivastava, “An efficient particle filtering technique on the Grassmann manifold,” in *Proceedings ICASSP*, Dallas, TX, March 14-19 2010, pp. 3838–3841.
- [20] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*. Princeton, NJ: Princeton University Press, 2008.
- [21] P.-A. Absil, “Optimization on manifolds: Methods and applications,” Université Catholique Louvain, Tech. Rep. UCL-INMA-2009.043, 2009.
- [22] T. E. Abrudan, J. Eriksson, and V. Koivunen, “Steepest descent algorithms for optimization under unitary matrix constraint,” *IEEE Transactions Signal Processing*, vol. 56, no. 3, pp. 1134–1147, March 2008.
- [23] —, “Conjugate gradient algorithm for optimization under unitary matrix constraint,” *Signal Processing*, vol. 89, pp. 1704–1714, 2009.
- [24] S. Fiori and T. Tanaka, “An algorithm to compute averages on matrix lie groups,” *IEEE Transactions Signal Processing*, vol. 57, no. 12, pp. 4734–4743, December 2009.
- [25] J. J. Boutros, F. Kharrat-Kammoun, and H. Randriambololona, “A classification of multiple antenna channels,” in *Proceedings International Zurich Seminar on Communications*, February 22-24 2006, pp. 14–17.
- [26] K. V. Mardia and P. E. Jupp, *Directional Statistics*. John Wiley & Sons, 1999.
- [27] Y. Chikuse, *Statistics on special manifolds*. New York: Springer Verlag, 2003.
- [28] K. Mammassis, R. W. Stewart, and J. S. Thompson, “Spatial fading correlation model using mixtures of von Mises Fisher distributions,” *IEEE Transactions Wireless Communications*, vol. 8, no. 4, pp. 2046–2055, April 2009.
- [29] K. Mammassis and R. W. Stewart, “The Fisher-Bingham spatial correlation model for multiple antenna systems,” *IEEE Transactions Vehicular Technology*, vol. 58, no. 5, pp. 2130–2136, June 2009.
- [30] O. Hamsici and A. Martinez, “Rotation invariant kernels and their application to shape analysis,” *IEEE Transactions Pattern Analysis Machine Intelligence*, vol. 31, no. 11, pp. 1985–1999, November 2009.
- [31] J. Ward, “Space-time adaptive processing for airborne radar,” Lincoln Laboratory, Massachusetts Institute of Technology, Lexington, MA, Tech. Rep. 1015, December 1994.
- [32] C. P. Robert and G. Casella, *Monte Carlo Statistical Methods*, 2nd ed. New York: Springer Verlag, 2004.
- [33] C. P. Robert, *The Bayesian Choice - From Decision-Theoretic Foundations to Computational Implementation*. New York: Springer Verlag, 2007.

- [34] P. D. Hoff, “Simulation of the matrix Bingham-von Mises-Fisher distribution, with applications to multivariate and relational data,” *Journal of Computational and Graphical Statistics*, vol. 18, no. 2, pp. 438–456, June 2009.
- [35] A. Sarlette and R. Sepulchre, “Consensus optimization on manifolds,” *SIAM Journal Control Optimization*, vol. 48, no. 1, pp. 56–76, 2009.
- [36] E. Begelfor and M. Werman, “Affine invariance revisited,” in *Proceedings IEEE CVPR06*, 2006, pp. 2087–2094.
- [37] N. Keshava and J. Mustard, “Spectral unmixing,” *IEEE Signal Processing Magazine*, vol. 19, no. 1, pp. 44–57, January 2002.
- [38] A. Halimi, Y. Altmann, N. Dobigeon, and J.-Y. Tourneret, “Nonlinear unmixing of hyperspectral images using a generalized bilinear model,” *IEEE Transactions Geoscience Remote Sensing*, 2011, to appear.
- [39] W. Fan, B. Hu, J. Miller, and M. Li, “Comparative study between a new nonlinear model and common linear model for analysing laboratory simulated-forest hyperspectral data,” *International Journal Remote Sensing*, vol. 30, no. 11, pp. 2951–2962, June 2009.
- [40] N. Dobigeon, S. Moussaoui, M. Coulon, J.-Y. Tourneret, and A. Hero, “Joint Bayesian endmember extraction and linear unmixing for hyperspectral imagery,” *IEEE Transactions on Signal Processing*, vol. 57, no. 11, pp. 4355–4368, November 2009.
- [41] G. Letac, “Familles exponentielles invariantes sur les matrices symétriques,” December 2010, private communication.
- [42] Y. Chung, “Simulation of truncated gamma variables,” *Korean J. Comput. Appl. Math.*, vol. 5, no. 3, pp. 601–610, 1998.
- [43] P. E. Jupp and K. V. Mardia, “Maximum likelihood estimators for the matrix von Mises-Fisher and Bingham distributions,” *The Annals of Statistics*, vol. 7, no. 3, pp. 599–606, May 1979.

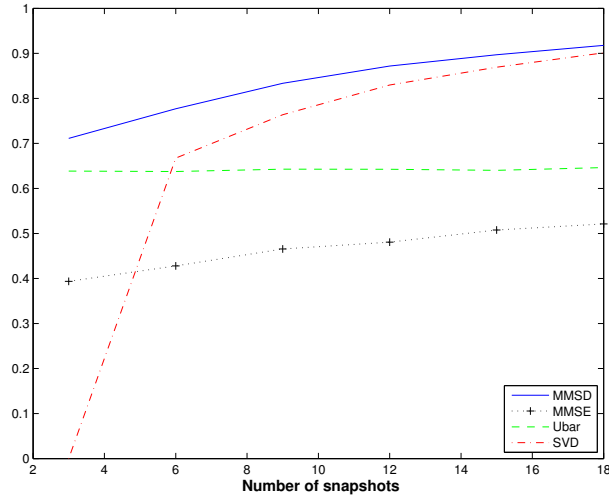


Figure 4: Fraction of energy of  $\hat{\mathbf{U}}$  in  $\mathcal{R}(\mathbf{U})$  versus  $K$ .  $N = 20$ ,  $p = 5$ ,  $\kappa = 20$  and  $SNR = 5\text{dB}$ . Linear model, Bingham prior.

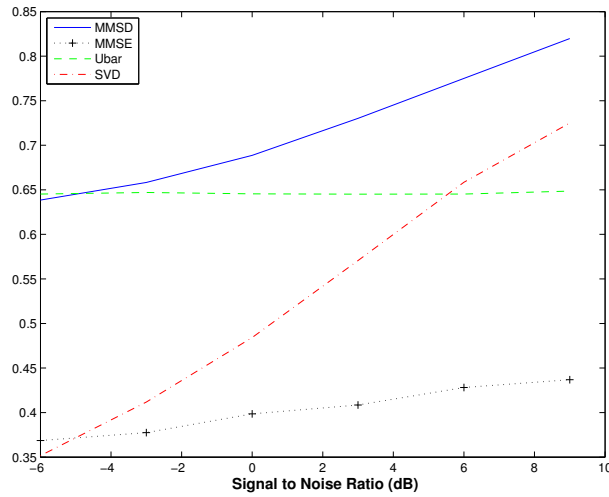


Figure 5: Fraction of energy of  $\hat{\mathbf{U}}$  in  $\mathcal{R}(\mathbf{U})$  versus  $SNR$ .  $N = 20$ ,  $p = 5$ ,  $\kappa = 20$  and  $K = 5$ . Linear model, Bingham prior.



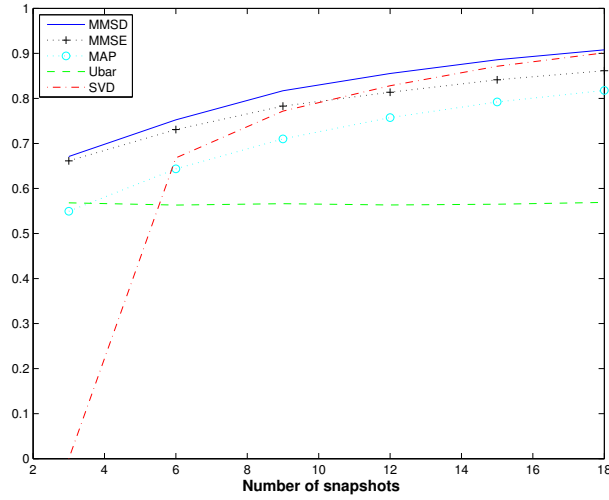


Figure 6: Fraction of energy of  $\hat{U}$  in  $\mathcal{R}(U)$  versus  $K$ .  $N = 20$ ,  $p = 5$ ,  $\kappa = 20$  and  $SNR = 5\text{dB}$ . Linear model, vMF prior.

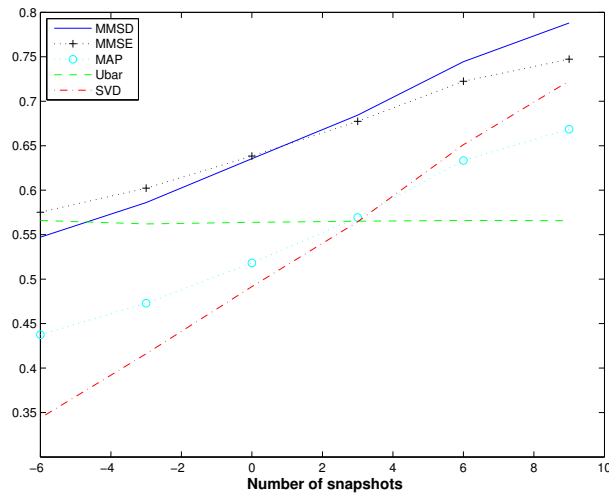


Figure 7: Fraction of energy of  $\hat{U}$  in  $\mathcal{R}(U)$  versus  $SNR$ .  $N = 20$ ,  $p = 5$ ,  $\kappa = 20$  and  $K = 5$ . Linear model, vMF prior.

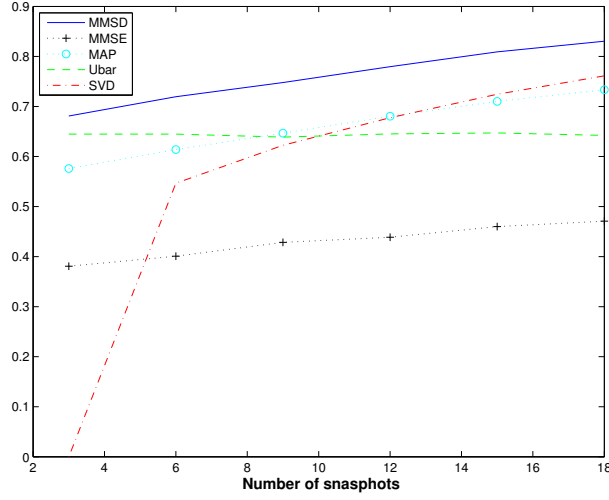


Figure 8: Fraction of energy of  $\hat{\mathbf{U}}$  in  $\mathcal{R}(\mathbf{U})$  versus  $K$ .  $N = 20$ ,  $p = 5$ ,  $\kappa = 20$ ,  $SNR_- = 5\text{dB}$  and  $SNR_+ = 10\text{dB}$ . Covariance matrix model, Bingham prior.

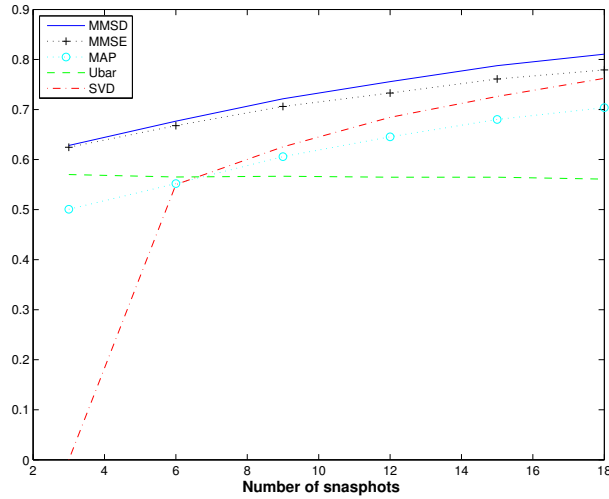


Figure 9: Fraction of energy of  $\hat{\mathbf{U}}$  in  $\mathcal{R}(\mathbf{U})$  versus  $K$ .  $N = 20$ ,  $p = 5$ ,  $\kappa = 20$ ,  $SNR_- = 5\text{dB}$  and  $SNR_+ = 10\text{dB}$ . Covariance matrix model, vMF prior.

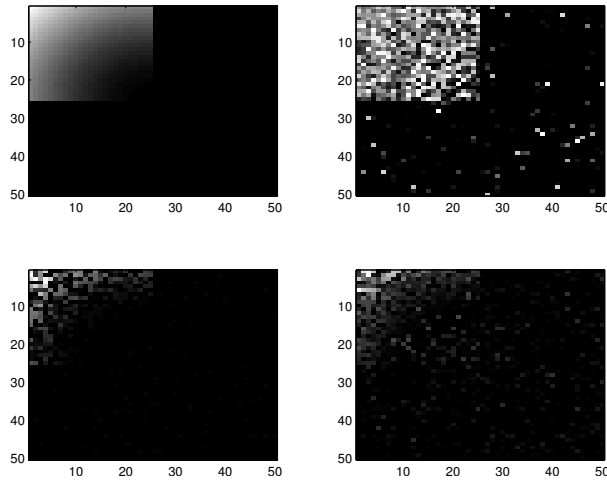


Figure 10: Top, left: non-linearity coefficients  $\gamma_{1,2}$ . Top, right: distance between  $\bar{U}$  and  $\hat{U}_n$  estimated with  $\eta = 0$ . Bottom: distance between  $\bar{U}$  and  $\hat{U}_\ell$  estimated with  $\eta = 0.5$  (left) and  $\eta = 50$  (right).

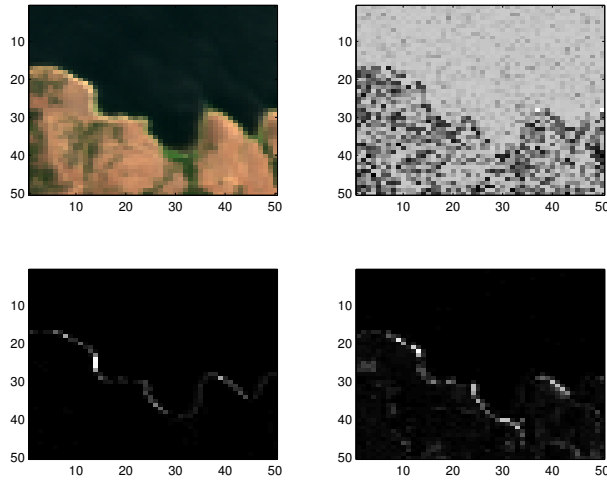


Figure 11: Top, left: The Moffett Field scene as composite true colors. Top, right: distance between  $\bar{U}$  and  $\hat{U}_n$  estimated with  $\eta = 0$ . Bottom: distance between  $\bar{U}$  and  $\hat{U}_\ell$  estimated with  $\eta = 0.5$  (left) and  $\eta = 50$  (right).