

# A Riemannian approach to blind separation of $t$ -distributed sources

Florent Bouchard\*, Arnaud Breloy<sup>†</sup>, Guillaume Ginolhac\*, Alexandre Renaux<sup>‡</sup>

\*LISTIC, Univ. Savoie Mont Blanc, France, <sup>†</sup>LEME, Univ. Paris Nanterre, France, <sup>‡</sup> L2S, Univ. Paris Sud, France

Email: florent.bouchard@univ-smb.fr

**Abstract**—The blind source separation problem is considered through the approach based on non-stationarity and coloration. In both cases, the sources are usually assumed to be Gaussian. In this paper, we extend previous works in order to handle sources drawn from the multivariate Student  $t$ -distribution. After studying the structure of the parameter manifold in this case, a new blind source separation criterion based on the log-likelihood of the considered distribution is proposed. To solve the resulting optimization problem, Riemannian optimization on the parameter manifold is leveraged. Practical expressions of the mathematical tools required by first order Riemannian optimization methods for this parameter manifold are derived to this end. The performance of the proposed method is illustrated on simulated data.

**Index Terms**—blind source separation, Student  $t$ -distribution, Riemannian optimization

## I. INTRODUCTION

Blind source separation (BSS) is an ubiquitous tool for signal processing and data analysis with applications in many engineering fields such as radar, communications, image processing and biomedical signal analysis; see [1] for a complete overview. In this work, we consider the determined linear instantaneous mixing model

$$\mathbf{x} = \mathbf{A}\mathbf{s}, \quad (1)$$

where  $\mathbf{x} \in \mathbb{R}^n$  corresponds to the observations,  $\mathbf{s} \in \mathbb{R}^n$  contains the sources and  $\mathbf{A} \in \text{GL}_n$  ( $n \times n$  non-singular matrices) is the mixing matrix. Given some observations of  $\mathbf{x}$ , the goal is to estimate the sources  $\mathbf{s}$  and the mixing process  $\mathbf{A}$ .

Up to now, two different types of methods have been proposed to perform blind source separation with model (1) depending on the data at hand. Historically, in the original formulation of the problem, which corresponds to independent component analysis (ICA), sources are assumed to be independent and identically distributed (i.i.d.) [2]–[4]. To ensure separability, sources have to be non-Gaussian and higher order statistics must be employed. On the other hand, if one wants to be able to separate Gaussian sources, the i.i.d. assumption must be removed, *i.e.*, some structure of the sources is exploited such as non-stationarity or coloration [5], [6].

In this paper, we are interested in the second approach, *i.e.*, when non-stationarity or coloration is exploited to perform blind source separation. In this case, one considers  $T$  realizations of  $K$  Gaussian random variables  $\mathbf{x}_k = \mathbf{A}\mathbf{s}_k$  with covariance matrices  $\mathbf{A}\mathbf{\Lambda}_k\mathbf{A}^T \in \mathcal{S}_n^{++}$  ( $n \times n$  symmetric positive definite matrices), where the  $\mathbf{s}_k$  are the random

variables corresponding to the sources, which are assumed to be uncorrelated and whose covariance matrices are  $\mathbf{\Lambda}_k \in \mathcal{D}_n^{++}$  ( $n \times n$  diagonal positive definite matrices). In practice, the  $K$  different sets of observations  $\{\mathbf{x}_k(t)\}$  can for instance correspond to different epochs (non-stationarity) or different frequencies (coloration) [1], [6], [7]. To solve the problem in the Gaussian case, one usually performs a two-steps procedure. First, the sample covariance matrices  $\mathbf{C}_k$  of the observations  $\{\mathbf{x}_k(t)\}$  are computed. Then, matrices  $\mathbf{C}_k$  are jointly diagonalized with a matrix  $\mathbf{B} \in \text{GL}_n$ , solution to the optimization problem

$$\underset{\mathbf{B} \in \text{GL}_n}{\text{argmin}} \sum_k f(\mathbf{B}, \mathbf{C}_k), \quad (2)$$

where  $f$  is a diagonality criterion of  $\mathbf{B}\mathbf{C}_k\mathbf{B}^T$ . Estimates  $\hat{\mathbf{s}}$  and  $\hat{\mathbf{A}}$  of sources  $\mathbf{s}$  and mixing matrix  $\mathbf{A}$  are finally given by  $\hat{\mathbf{s}} = \mathbf{B}\mathbf{x}$  and  $\hat{\mathbf{A}} = \mathbf{B}^{-1}$ , up to the usual permutation and diagonal scaling ambiguities.

Previous research on this topic have mainly focused on two aspects: the choice of the diagonality criterion  $f$  in (2) and the conception of iterative methods to solve (2) given a criterion  $f$ . Concerning the diagonality criterion  $f$ , most articles employ either the least squares one proposed in [4] or the one based on the log-likelihood introduced in [6]. Recently, different criteria based on several divergences and distances on  $\mathcal{S}_n^{++}$  have been proposed in [8], [9]. Many algorithms to solve (2) for various diagonality criteria  $f$  have been developed; see *e.g.*, [4], [6], [9]–[12]. In particular, methods based on Riemannian optimization can be found in [9], [10], [12].

Works based on non-stationarity and coloration consider that sources are Gaussian. In this paper, our main contribution is to extend this approach to handle sources drawn from the multivariate Student  $t$ -distribution [13]. Since it has a heavier tail than the Gaussian distribution, it is more adapted to account for outliers. To do so, we derive a new blind source separation criterion based on the log-likelihood of the data model. To find a minimizer of this criterion, we propose to exploit Riemannian optimization on the parameter manifold, which allows to find estimates  $\hat{\mathbf{A}} \in \text{GL}_n$  and  $\{\hat{\mathbf{\Lambda}}_k\} \in (\mathcal{D}_n^{++})^K$  of the true parameters  $\mathbf{A}$  and  $\{\mathbf{\Lambda}_k\}$ . In this scope, we derive practical expressions for the mathematical tools required by first order based Riemannian optimization methods on this parameter manifold.

This paper contains four sections including this introduction. In section II, the data model is defined along with the

manifold holding the parameters of the considered distribution. To handle the diagonal scaling ambiguity of blind source separation, an original constraint is proposed. Finally, the Riemannian geometry of the resulting manifold is studied. In section III, a Riemannian optimization framework is developed in order to be able to solve optimization problems on the parameter manifold of interest. A blind source separation criterion based on the log-likelihood function of the considered distribution is then derived, thus yielding an original blind source separation method. In section IV, the performance of the proposed method is illustrated on simulated data. As expected, when sources are drawn from the multivariate Student  $t$ -distribution, the proposed method outperforms the state of the art method [6].

## II. MODEL

To design a new blind source separation method adapted to the multivariate Student  $t$ -distribution, we first need to clearly establish the model of the data along with the parameters they depend on. In section II-A, the probability density function of the distribution of the data and the manifold holding the parameters of interest are defined. In section II-B, as the approach chosen to solve the blind source separation problem is based on Riemannian optimization, the Riemannian geometry of the parameter manifold is studied.

### A. Data distribution and parameter manifold

We consider  $K$  sets of  $T$  observations  $\{\mathbf{x}_k(t)\}$  of the random variables  $\mathbf{x}_k$  in  $\mathbb{R}^n$  following the elliptical distributions [13] with density generator function  $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  and covariance matrices  $\mathbf{A}\mathbf{\Lambda}_k\mathbf{A}^T$ , where  $\mathbf{A} \in \text{GL}_n$  and  $\mathbf{\Lambda}_k \in \mathcal{D}_n^{++}$ . The probability density function  $f$  of  $\{\mathbf{x}_k(t)\}$  is

$$f(\{\mathbf{x}_k(t)\}|\mathbf{A}, \{\mathbf{\Lambda}_k\}) = \prod_k f_d(\{\mathbf{x}_k(t)\}|\mathbf{A}\mathbf{\Lambda}_k\mathbf{A}^T), \quad (3)$$

where  $f_d$  is the probability density function of the multivariate Student  $t$ -distribution with  $d \in \mathbb{N}^*$  degrees of freedom. Given  $\{\mathbf{x}(t)\}$  with covariance matrix  $\mathbf{C}$ , we have, up to a factor,

$$f_d(\{\mathbf{x}(t)\}|\mathbf{C}) = \prod_t \det(\mathbf{C})^{-1/2} \left(1 + \frac{1}{d} \mathbf{x}(t)^T \mathbf{C}^{-1} \mathbf{x}(t)\right)^{-(d+n)/2}. \quad (4)$$

The set of parameters  $(\mathbf{A}, \mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_K)$  of the random variables  $\mathbf{x}_k$  lies in the manifold  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$ . The parameters  $(\mathbf{A}, \mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_K)$  of the  $\mathbf{x}_k$  are not unique: given any permutation matrix  $\mathbf{P}$  and non-singular diagonal matrix  $\mathbf{\Sigma}$ , we have

$$f(\{\mathbf{x}_k(t)\}|\mathbf{A}, \{\mathbf{\Lambda}_k\}) = f(\{\mathbf{x}_k(t)\}|\mathbf{A}\mathbf{P}\mathbf{\Sigma}, \{\mathbf{\Sigma}^{-1}\mathbf{P}^T\mathbf{\Lambda}_k\mathbf{P}\mathbf{\Sigma}^{-1}\}).$$

To deal with the diagonal scaling ambiguity resulting from the invariance through the action of  $\mathbf{\Sigma}$ , an additional constraint can

be imposed. In this work, we propose a new constraint, related to the one introduced in [11]. It is given by

$$\sum_k \mathbf{\Lambda}_k = \mathbf{I}_n, \quad (5)$$

where  $\mathbf{I}_n$  is the identity matrix. As this constraint is smooth with respect to  $(\mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_K)$ , it can be integrated in a smooth manifold. The permutation ambiguity resulting from the invariance through the action of  $\mathbf{P}$  is more tricky as it cannot be handled in a manifold. However, contrary to the diagonal scaling ambiguity, it is not an issue from a numerical point of view as it does not yield algorithms to converge to degenerate solutions. As most works on blind source separation, we choose to ignore it. It follows that the manifold holding the parameters of the distribution of interest is chosen as

$$\mathcal{M} = \{(\mathbf{A}, \{\mathbf{\Lambda}_k\}) \in \text{GL}_n \times (\mathcal{D}_n^{++})^K : \sum_k \mathbf{\Lambda}_k = \mathbf{I}_n\}, \quad (6)$$

which is a submanifold of  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$ .

### B. Riemannian geometry of the parameter manifold

In the following,  $\theta = (\mathbf{A}, \{\mathbf{\Lambda}_k\})$ ,  $\xi = (\xi_{\mathbf{A}}, \{\xi_k\})$  and  $\eta = (\eta_{\mathbf{A}}, \{\eta_k\})$ . Since  $\text{GL}_n$  and  $\mathcal{D}_n^{++}$  are open in  $\mathbb{R}^{n \times n}$  and  $\mathcal{D}_n$  ( $n \times n$  diagonal matrices), respectively, the tangent space  $T_\theta \text{GL}_n \times (\mathcal{D}_n^{++})^K$  can be identified as  $\mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$ . From a statistical point of view, one ideally wants to choose the Fisher information metric associated with (3) as the Riemannian metric of  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$ . However, as illustrated in [14] for Gaussian sources, the Fisher information metric appears complicated to deal with in this case and a different Riemannian metric is preferred. In this work, we choose

$$\langle \xi, \eta \rangle_\theta = \text{tr}((\xi_{\mathbf{A}} \mathbf{A}^{-1})^T \eta_{\mathbf{A}} \mathbf{A}^{-1}) + \sum_k \text{tr}(\mathbf{\Lambda}_k^{-2} \xi_k \eta_k). \quad (7)$$

The component that depends on  $\mathbf{A}$  corresponds to the so-called right-invariant metric on  $\text{GL}_n$  [12], [15]. The component that depends on  $\mathbf{\Lambda}_k$  corresponds to the popular affine-invariant metric on  $\mathcal{D}_n^{++}$  [16]. The main interest of the chosen metric on  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$  is that it is invariant with respect to the actions of diagonal and permutation matrices  $\mathbf{\Sigma}$  and  $\mathbf{P}$  described in the previous section.

The manifold  $\mathcal{M}$  defined in (6) that holds the parameters of interest is a submanifold of  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$ . Thus, the tangent space  $T_\theta \mathcal{M}$  at  $\theta \in \mathcal{M}$  is a subspace of  $T_\theta \text{GL}_n \times (\mathcal{D}_n^{++})^K$ . Furthermore,  $\mathcal{M}$  can inherit the Riemannian metric (7) of  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$ . Given  $\theta \in \mathcal{M}$ , the tangent space  $T_\theta \mathcal{M}$  along with the orthogonal projection map  $P_\theta : \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K \rightarrow T_\theta \mathcal{M}$  according to (7) are given in Proposition 1.

**Proposition 1.** *The tangent space  $T_\theta \mathcal{M}$  at  $\theta \in \mathcal{M}$  is*

$$T_\theta \mathcal{M} = \{\xi \in \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K : \sum_k \xi_k = \mathbf{0}\}.$$

*The orthogonal projection map  $P_\theta : \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K \rightarrow T_\theta \mathcal{M}$  according to (7) is given by*

$$P_\theta(\xi) = (\xi_{\mathbf{A}}, \{\xi_k - (\sum_\ell \mathbf{\Lambda}_\ell^2)^{-1} \sum_\ell \xi_\ell \mathbf{\Lambda}_\ell^2\}).$$

*Proof.* The manifold  $\mathcal{M}$  of interest is defined as  $\mathcal{M} = \{\theta \in \text{GL}_n \times (\mathcal{D}_n^{++})^K : \varphi(\theta) = \mathbf{0}\}$ , where  $\varphi : \text{GL}_n \times (\mathcal{D}_n^{++})^K \rightarrow \mathbb{R}^{n \times n}$  is the smooth mapping such that  $\varphi(\theta) = \sum_k \mathbf{\Lambda}_k - \mathbf{I}_n$ . It follows that the tangent space  $T_\theta \mathcal{M}$  at  $\theta \in \mathcal{M}$  is  $T_\theta \mathcal{M} = \{\xi \in \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K : \text{D}\varphi(\theta)[\xi] = \mathbf{0}\}$ , where  $\text{D}\varphi(\theta)[\xi]$  is the derivative of  $\varphi$  at  $\theta$  in direction  $\xi$ . Noticing that  $\text{D}\varphi(\theta)[\xi] = \sum_k \xi_k$  yields the result.

The normal space  $N_\theta \mathcal{M}$  at  $\theta \in \mathcal{M}$  is defined as  $N_\theta \mathcal{M} = \{\eta \in \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K : \forall \xi \in T_\theta \mathcal{M}, \langle \xi, \eta \rangle_\theta = 0\}$ . One can check that, in our case, it is given by

$$N_\theta \mathcal{M} = \{(\mathbf{0}, \{\mathbf{\Lambda}_k^2 \boldsymbol{\eta}\}) : \boldsymbol{\eta} \in \mathcal{D}_n\}.$$

It follows that  $P_\theta(\xi) = \xi - (\mathbf{0}, \{\mathbf{\Lambda}_k^2 \boldsymbol{\eta}\})$ , where  $\boldsymbol{\eta}$  is the diagonal matrix such that  $P_\theta(\xi) \in T_\theta \mathcal{M}$ , i.e.,  $\sum_k (\xi_k - \mathbf{\Lambda}_k^2 \boldsymbol{\eta}) = \mathbf{0}$ . We thus obtain  $\boldsymbol{\eta} = (\sum_\ell \mathbf{\Lambda}_\ell^2)^{-1} \sum_\ell \xi_\ell$ . ■

These elements of Riemannian geometry are sufficient to build the Riemannian optimization framework on  $\mathcal{M}$  required for blind source separation in section III.

### III. BLIND SOURCE SEPARATION METHOD

The blind source separation problem can be written as an optimization problem on the manifold  $\mathcal{M}$ . In section III-A, a Riemannian optimization framework on  $\mathcal{M}$  is built. In section III-B, a blind source separation based on the likelihood of the multivariate Student  $t$ -distribution is proposed and the elements required for Riemannian optimization on  $\mathcal{M}$  are provided.

#### A. Riemannian optimization on the parameter manifold

Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be an objective function defined on  $\mathcal{M}$ . To be able to perform Riemannian optimization [17] of  $f$  on  $\mathcal{M}$ , it is needed to have the Riemannian gradient of  $f$  along with a retraction on  $\mathcal{M}$ , which, given  $\theta \in \mathcal{M}$ , maps the elements of  $T_\theta \mathcal{M}$  back onto  $\mathcal{M}$ . Concerning the Riemannian gradient on  $\mathcal{M}$ , Proposition 2 shows that it can be obtained from the Euclidean gradient of  $f$ , which is defined on the ambient space  $\mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$ .

**Proposition 2.** *The Riemannian gradient  $\text{grad}_{\mathcal{M}} f(\theta) \in T_\theta \mathcal{M}$  of  $f$  at  $\theta \in \mathcal{M}$  according to metric (7) is given by*

$$\text{grad}_{\mathcal{M}} f(\theta) = P_\theta \left( \nabla f_{\mathbf{A}} \mathbf{A}^T \mathbf{A}, \{\mathbf{\Lambda}_k^2 \nabla f_k\} \right),$$

where  $P_\theta : \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K \rightarrow T_\theta \mathcal{M}$  is defined in Proposition 1 and  $\nabla f(\theta) = (\nabla f_{\mathbf{A}}, \{\nabla f_k\}) \in \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$  is the Euclidean gradient of  $f$ .

*Proof.* First recall that  $\nabla f(\theta)$  is the only element of  $\mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$  such that for all  $\xi \in \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$ ,

$$\text{D}f(\theta)[\xi] = \langle \nabla f(\theta), \xi \rangle^\mathcal{E},$$

where  $\langle \cdot, \cdot \rangle^\mathcal{E}$  is the Euclidean metric on  $\mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$ , which is defined as

$$\langle \xi, \eta \rangle^\mathcal{E} = \text{tr}(\boldsymbol{\xi}_\mathbf{A}^T \boldsymbol{\eta}_\mathbf{A}) + \sum_k \text{tr}(\boldsymbol{\xi}_k \boldsymbol{\eta}_k).$$

Similarly, the Riemannian gradient of  $f$  on  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$  according to the Riemannian metric (7) is the only element of  $\mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$  such that, for all  $\xi \in \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$ ,

$$\text{D}f(\theta)[\xi] = \langle \text{grad}_{\text{GL}_n \times (\mathcal{D}_n^{++})^K} f(\theta), \xi \rangle_\theta,$$

From  $\langle \text{grad}_{\text{GL}_n \times (\mathcal{D}_n^{++})^K} f(\theta), \xi \rangle_\theta = \langle \nabla f(\theta), \xi \rangle^\mathcal{E}$ , one can check that

$$\text{grad}_{\text{GL}_n \times (\mathcal{D}_n^{++})^K} f(\theta) = \left( \nabla f_{\mathbf{A}} \mathbf{A}^T \mathbf{A}, \{\mathbf{\Lambda}_k^2 \nabla f_k\} \right).$$

The rest of the proof follows from [17, Equation (3.37)], which indicates that  $\text{grad}_{\mathcal{M}} f(\theta)$  is obtained by projecting  $\text{grad}_{\text{GL}_n \times (\mathcal{D}_n^{++})^K} f(\theta)$  onto  $T_\theta \mathcal{M}$ . ■

The Riemannian gradient  $\text{grad}_{\mathcal{M}} f(\theta) \in T_\theta \mathcal{M}$  of  $f$  at  $\theta \in \mathcal{M}$  allows to define a descent direction  $\xi$  of  $f$  in the tangent space  $T_\theta \mathcal{M}$ . In order to achieve a new point on the manifold  $\mathcal{M}$  from  $\xi$ , a retraction  $R$  on  $\mathcal{M}$  is needed. Given  $\theta \in \mathcal{M}$ , the retraction  $R_\theta$  at  $\theta$  is the mapping from  $T_\theta \mathcal{M}$  onto  $\mathcal{M}$  such that:

- 1)  $R_\theta(0_\theta) = \theta$ , where  $0_\theta$  is the zero element of  $T_\theta \mathcal{M}$ .
- 2)  $\text{D}R_\theta(0_\theta)[\xi] = \xi$ .

Proposition 3 contains a proper retraction on  $\mathcal{M}$  obtained by projecting a second order approximation of the Riemannian exponential map of  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$  onto  $\mathcal{M}$ .

**Proposition 3.** *Given  $\theta \in \mathcal{M}$ , a proper retraction  $R_\theta : T_\theta \mathcal{M} \rightarrow \mathcal{M}$  at  $\theta$  is given, for all  $\xi \in T_\theta \mathcal{M}$  by*

$$R_\theta(\xi) = \Gamma(\bar{R}_\theta(\xi)),$$

where  $\bar{R}_\theta$  is a second order approximation of the Riemannian exponential map of  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$  equipped with (7), which is given by

$$\bar{R}_\theta(\xi) = \theta + \xi +$$

$$\frac{1}{2} (\boldsymbol{\xi}_\mathbf{A} \mathbf{A}^{-1} \boldsymbol{\xi}_\mathbf{A} + \boldsymbol{\xi}_\mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \boldsymbol{\xi}_\mathbf{A}^T \mathbf{A} - \mathbf{A}^{-T} \boldsymbol{\xi}_\mathbf{A}^T \boldsymbol{\xi}_\mathbf{A}, \{\mathbf{\Lambda}_k^{-1} \xi_k^2\}),$$

and  $\Gamma : \text{GL}_n \times (\mathcal{D}_n^{++})^K \rightarrow \mathcal{M}$  is the projection map defined, for all  $\theta \in \text{GL}_n \times (\mathcal{D}_n^{++})^K$ , as

$$\Gamma(\theta) = (\mathbf{A}, \{(\sum_\ell \Lambda_\ell)^{-1} \Lambda_k\}).$$

*Proof.* Since  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$  equipped with metric (7) is a Riemannian product manifold, we know from [17] that the Riemannian exponential map is defined, for all  $\theta \in \text{GL}_n \times (\mathcal{D}_n^{++})^K$  and  $\xi \in \mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$ , as

$$\exp_{\text{GL}_n \times (\mathcal{D}_n^{++})^K}(\xi) = \left( \exp_{\mathbf{A}}^{\text{GL}_n}(\boldsymbol{\xi}_\mathbf{A}), \{\exp_{\mathbf{\Lambda}_k}^{\mathcal{D}_n^{++}}(\boldsymbol{\xi}_k)\} \right),$$

where

$$\exp_{\mathbf{A}}^{\text{GL}_n}(\boldsymbol{\xi}_\mathbf{A}) =$$

$$\exp(\boldsymbol{\xi}_\mathbf{A} \mathbf{A}^{-1} - (\boldsymbol{\xi}_\mathbf{A} \mathbf{A}^{-1})^T) \exp(\boldsymbol{\xi}_\mathbf{A} \mathbf{A}^{-1})^T \mathbf{A}$$

is the Riemannian exponential map of  $\text{GL}_n$  equipped with the right-invariant metric [12], [15] and

$$\exp_{\mathbf{\Lambda}_k}^{\mathcal{D}_n^{++}}(\boldsymbol{\xi}_k) = \mathbf{\Lambda}_k \exp(\mathbf{\Lambda}_k^{-1} \boldsymbol{\xi}_k)$$

is the Riemannian exponential map of  $\mathcal{D}_n^{++}$  equipped with the affine-invariant metric [16]. The matrix exponential admits the

second order approximation

$$\exp(\mathbf{X}) = \mathbf{I}_n + \mathbf{X} + \frac{1}{2}\mathbf{X}^2 + o(\mathbf{X}^2).$$

By injecting this approximation in the formula of the Riemannian exponential map on  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$ , one can show that the proposed operator  $\bar{R}$  is a second order approximation of this Riemannian exponential map, which is also enough to prove that it is a retraction on  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$ . Moreover, it is readily checked that  $\Gamma$  is a projection map from  $\text{GL}_n \times (\mathcal{D}_n^{++})^K$  onto  $\mathcal{M}$ . Therefore, from [18], we know that the proposed operator  $R$  is a proper retraction on  $\mathcal{M}$ . ■

The Riemannian gradient of  $f : \mathcal{M} \rightarrow \mathbb{R}$  and the retraction  $R$  on  $\mathcal{M}$  are enough to define a Riemannian gradient descent algorithm on  $\mathcal{M}$ . Indeed, given iterate  $\theta_i$ , the next iterate is obtained by

$$\theta_{i+1} = R_{\theta_i}(-t_i \text{grad}_{\mathcal{M}} f(\theta_i)), \quad (8)$$

where  $t_i$  is the stepsize, which can for instance be computed through a linesearch; see [17] for details.

To be able to use more sophisticated Riemannian optimization algorithms such as conjugate gradient or BFGS, one further needs to define a vector transport operator  $\mathcal{T}$ , which allows to transport a tangent vector from one tangent space onto another [17]. Given  $\theta \in \mathcal{M}$  and  $\xi, \eta \in T_{\theta}\mathcal{M}$ , the vector transport  $\mathcal{T}(\theta, \xi, \eta)$  associated to the retraction  $R$  transports the tangent vector  $\eta$  in the tangent space  $T_{R_{\theta}(\xi)}\mathcal{M}$  of  $R_{\theta}(\xi)$ . A generic solution [17] is to choose  $\mathcal{T}(\theta, \xi, \eta)$  as

$$\mathcal{T}(\theta, \xi, \eta) = P_{R_{\theta}(\xi)}(\eta), \quad (9)$$

where, in our case,  $P$  is the projection map defined in Proposition 1 and  $R$  is the retraction defined in Proposition 3.

### B. Application to the log-likelihood of the $t$ -distribution

To perform the blind source separation of the  $K$  sets of  $T$  observations  $\{\mathbf{x}_k(t)\}$  drawn from the distribution defined in Section II-A, the optimal solution is to employ the maximum likelihood estimator associated with the probability density function (3). It is equivalent to minimizing the negative log-likelihood of (3) on  $\mathcal{M}$ , which is given by

$$f(\theta) = \sum_k f_{\mathbf{x}_k}^{++}(\psi_k(\theta)), \quad (10)$$

where  $\psi_k : \mathcal{M} \rightarrow \mathcal{S}_n^{++}$  is such that  $\psi_k(\theta) = \mathbf{A}\mathbf{\Lambda}_k\mathbf{A}^T$  and  $f^{++}$  is the negative log-likelihood of the multivariate Student  $t$ -distribution with  $d \in \mathbb{N}^*$  degrees of freedom. Given  $T$  observations  $\{\mathbf{x}(t)\}$ , the negative log-likelihood  $f_{\mathbf{x}}^{++} : \mathcal{S}_n^{++} \rightarrow \mathbb{R}$  associated with (4) is defined, up to constants, as

$$f_{\mathbf{x}}^{++}(\mathbf{C}) = \frac{T}{2} \log \det(\mathbf{C}) + \frac{d+n}{2} \sum_t \log(d + \mathbf{x}(t)^T \mathbf{C}^{-1} \mathbf{x}(t)). \quad (11)$$

To be able to minimize (10) with a Riemannian optimization algorithm on  $\mathcal{M}$  within the framework developed in

section III-A, it remains to provide the Euclidean gradient of  $f$  on  $\mathbb{R}^{n \times n} \times (\mathcal{D}_n)^K$ . This is achieved in Proposition 4.

**Proposition 4.** *The Euclidean gradient  $\nabla f(\theta)$  of the blind source separation criterion  $f$  defined in (10) at  $\theta \in \mathcal{M}$  is given by*

$$\nabla f(\theta) = \left( 2 \sum_k \nabla f_{\mathbf{x}_k}^{++}(\psi_k(\theta)) \mathbf{A} \mathbf{\Lambda}_k, \{ \text{ddiag}(\mathbf{A}^T \nabla f_{\mathbf{x}_k}^{++}(\psi_k(\theta)) \mathbf{A}) \} \right),$$

where  $\text{ddiag}(\cdot)$  returns the diagonal part of its argument and  $\nabla f_{\mathbf{x}}^{++}(\mathbf{C})$  is the Euclidean gradient of (11) at  $\mathbf{C} \in \mathcal{S}_n^{++}$ , which is

$$\nabla f_{\mathbf{x}}^{++}(\mathbf{C}) = \mathbf{C}^{-1} \left[ \frac{T}{2} \mathbf{C} - \frac{d+n}{2} \sum_t \frac{\mathbf{x}(t) \mathbf{x}(t)^T}{d + \mathbf{x}(t)^T \mathbf{C}^{-1} \mathbf{x}(t)} \right] \mathbf{C}^{-1}.$$

*Proof.* By definition of  $f$ ,

$$\begin{aligned} D f(\theta)[\xi] &= \sum_k D f_{\mathbf{x}_k}(\psi_k(\theta)) [D \psi_k(\theta)[\xi]] \\ &= \sum_k \langle \nabla f_{\mathbf{x}_k}^{++}(\psi_k(\theta)), D \psi_k(\theta)[\xi] \rangle^{\mathbb{R}^{n \times n}}, \end{aligned}$$

where  $D \psi_k(\theta)[\xi] = \mathbf{A} \mathbf{\Lambda}_k \xi_{\mathbf{A}}^T + \xi_{\mathbf{A}} \mathbf{\Lambda}_k \mathbf{A}^T + \mathbf{A} \xi_k \mathbf{A}^T$  and  $\langle \cdot, \cdot \rangle^{\mathbb{R}^{n \times n}}$  is the Euclidean inner product on  $\mathbb{R}^{n \times n}$ . Basic calculations yield

$$\begin{aligned} D f(\theta)[\xi] &= \sum_k 2 \langle \nabla f_{\mathbf{x}_k}^{++}(\psi_k(\theta)), \xi_{\mathbf{A}} \mathbf{\Lambda}_k \mathbf{A}^T \rangle^{\mathbb{R}^{n \times n}} \\ &\quad + \sum_k \langle \nabla f_{\mathbf{x}_k}^{++}(\psi_k(\theta)), \mathbf{A} \xi_k \mathbf{A}^T \rangle^{\mathbb{R}^{n \times n}} \\ &= \langle 2 \sum_k \nabla f_{\mathbf{x}_k}^{++}(\psi_k(\theta)) \mathbf{A} \mathbf{\Lambda}_k, \xi_{\mathbf{A}} \rangle^{\mathbb{R}^{n \times n}} \\ &\quad + \sum_k \langle \text{ddiag}(\mathbf{A}^T \nabla f_{\mathbf{x}_k}^{++}(\psi_k(\theta)) \mathbf{A}), \xi_k \rangle^{\mathbb{R}^{n \times n}}, \end{aligned}$$

where we used  $\text{tr}(\mathbf{X}\mathbf{\Sigma}) = \text{tr}(\text{ddiag}(\mathbf{X})\mathbf{\Sigma})$  for all  $\mathbf{\Sigma} \in \mathcal{D}_n$ . The result is finally obtained by identification. ■

## IV. NUMERICAL ILLUSTRATION

To simulate data, we generate  $K = 30$   $n \times n$  (with  $n = 10$ ) symmetric positive definite matrices  $\mathbf{C}_k$  according to model

$$\mathbf{C}_k = \mathbf{A} \mathbf{\Lambda}_k \mathbf{A}^T, \quad (12)$$

where  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , with  $\mathbf{U}$  and  $\mathbf{V}$  random orthogonal matrices, and  $\mathbf{\Sigma}$  random diagonal matrix whose minimal and maximal elements are  $1/\sqrt{\alpha}$  and  $\sqrt{\alpha}$ , where  $\alpha = 10$  is the condition number of  $\mathbf{A}$ . Matrices  $\mathbf{\Lambda}_k \in \mathcal{D}_n^{++}$  correspond to the covariance matrices of the sources and have i.i.d. elements drawn from the chi squared distribution with expectation one. For  $T \in \{15, 25, 50, 75, 100, 500, 1000\}$ , 100 sets of  $T$  realizations  $\{\mathbf{x}_k(t)\}$  are randomly drawn from the multivariate Student  $t$ -distribution with  $d = 3$  degrees of freedom.

In each case, the blind source separation of the simulated data is performed with two methods. For the first one, denoted *jadiag*, the sample covariance matrices  $\hat{\mathbf{C}}_k = T^{-1} \sum_t \mathbf{x}_k(t) \mathbf{x}_k(t)^T$  are computed. These matrices are then jointly diagonalized with  $\mathbf{B}$ , minimizer of the log-likelihood criterion for Gaussian sources [6]. The estimated mixing matrix and source covariance matrices are given by  $\hat{\mathbf{A}} = \mathbf{B}^{-1}$  and  $\hat{\mathbf{\Lambda}}_k = \text{ddiag}(\mathbf{B} \hat{\mathbf{C}}_k \mathbf{B}^T)$ . The second method, denoted

*MLE t-dist*, consists in optimizing the negative of the log-likelihood (10) on  $\mathcal{M}$  for  $d = 3$  with a Riemannian conjugate gradient algorithm [17]. Optimization is performed with manopt toolbox [19].

To measure the performance of the two methods, three performance indexes are considered. The first one is the popular Moreau-Macchi index [20], which measures how close  $M = \hat{A}^{-1}A$  is to a permuted diagonal matrix. It is defined as

$$I_{M-\Lambda}(M) = \frac{1}{2n(n-1)} \sum_{p=1}^n \left( \frac{\sum_{q=1}^n |M_{pq}|}{\max_{1 \leq q \leq n} |M_{pq}|} - 1 \right) + \frac{1}{2n(n-1)} \sum_{p=1}^n \left( \frac{\sum_{q=1}^n |M_{qp}|}{\max_{1 \leq q \leq n} |M_{qp}|} - 1 \right). \quad (13)$$

It is a measure between zero and one, with zero indicating a perfect recovery of the mixing process. The second one measures the distances in  $\mathcal{S}_n^{++}$  between  $A\Lambda_k A^T$  and  $\hat{A}\hat{\Lambda}_k \hat{A}^T$ . It is defined as

$$I_{(\mathcal{S}_n^{++})^\kappa} = \frac{1}{K} \sum_k \delta_R^2(A\Lambda_k A^T, \hat{A}\hat{\Lambda}_k \hat{A}^T), \quad (14)$$

where  $\delta_R$  is the usual Riemannian distance function on  $\mathcal{S}_n^{++}$  [16]. The last performance index, denoted  $I_{(s\mathcal{S}_n^{++})^\kappa}$ , is the same as  $I_{(\mathcal{S}_n^{++})^\kappa}$  except that we normalize the determinants of  $A\Lambda_k A^T$  and  $\hat{A}\hat{\Lambda}_k \hat{A}^T$ . It allows to only consider the structure of the estimated matrices without taking into account their power.

In Figure 1, we observe that the proposed method *MLE t-dist* features better performance than *jadiag* for all three performance indexes for every  $T$ . These results are expected since *MLE t-dist* corresponds to the optimal log-likelihood estimator according to these simulated data. Notice that the index  $I_{(\mathcal{S}_n^{++})^\kappa}$  of *jadiag* shows poor performance whereas  $I_{(s\mathcal{S}_n^{++})^\kappa}$  is more satisfying. It indicates that *jadiag* fails at retrieving the power of  $A\Lambda_k A^T$  but captures their structure. On the other hand, *MLE t-dist* estimates both the power and the structure of  $A\Lambda_k A^T$  for these simulated data.

## REFERENCES

- [1] P. Comon and C. Jutten. *Handbook of Blind Source Separation: Independent Component Analysis and Applications*. Academic Press, 1st edition, 2010.
- [2] C. Jutten and J. Herault. Blind separation of sources, part I: An adaptive algorithm based on neuromimetic architecture. *Signal processing*, 24(1):1–10, 1991.
- [3] P. Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994.
- [4] J.-F. Cardoso and A. Souloumiac. Blind beamforming for non Gaussian signals. *IEEE Proceedings-F*, 140(6):362–370, dec 1993.
- [5] A. Belouchrani, K. Abed-Meraim, J.-F. Cardoso, and E. Moulines. A blind source separation technique using second-order statistics. *IEEE Transactions on signal processing*, 45(2):434–444, 1997.
- [6] D. T. Pham and J.-F. Cardoso. Blind separation of instantaneous mixtures of nonstationary sources. *IEEE Transactions on Signal Processing*, 49(9):1837–1848, 2001.

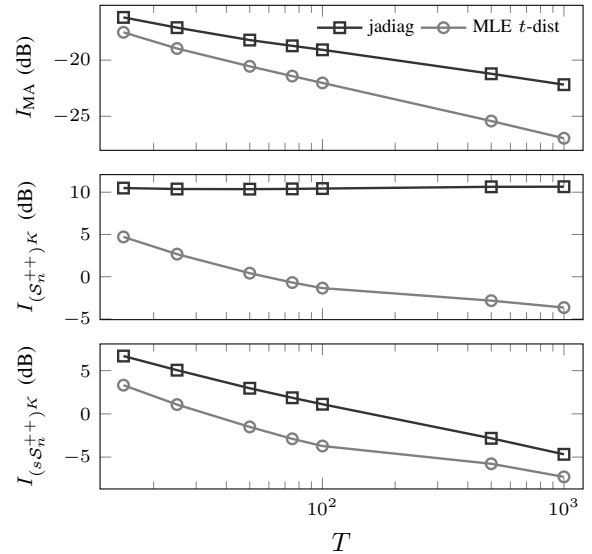


Fig. 1. Mean over 100 repetitions of the three performance indexes of the two considered blind source separation methods as functions of the number of samples  $T$ .

- [7] M. Congedo, C. Gouy-Pailler, and C. Jutten. On the blind source separation of human electroencephalogram by approximate joint diagonalization of second order statistics. *Clinical Neurophysiology*, 119(12):2677–2686, 2008.
- [8] K. Alyani, M. Congedo, and M. Moakher. Diagonality measures of Hermitian positive-definite matrices with application to the approximate joint diagonalization problem. *Linear Algebra and its Applications*, 2016.
- [9] F. Bouchard, J. Malick, and M. Congedo. Riemannian optimization and approximate joint diagonalization for blind source separation. *IEEE Transactions on Signal Processing*, 66(8):2041–2054, 2018.
- [10] B. Afsari. Sensitivity analysis for the problem of matrix joint diagonalization. *SIAM Journal on Matrix Analysis and Applications*, 30(3):1148–1171, 2008.
- [11] D.-T. Pham and M. Congedo. Least square joint diagonalization of matrices under an intrinsic scale constraint. In *Independent Component Analysis and Signal Separation*, pages 298–305. Springer, 2009.
- [12] F. Bouchard, B. Afsari, J. Malick, and M. Congedo. Approximate joint diagonalization with Riemannian optimization on the general linear group. *SIAM Journal of Matrix Analysis and Applications*, 41(1):152–170, 2020.
- [13] E. Ollila, D. E. Tyler, V. Koivunen, and H. V. Poor. Complex elliptically symmetric distributions: Survey, new results and applications. *IEEE Transactions on Signal Processing*, 60(11):5597–5625, 2012.
- [14] F. Bouchard, A. Breloy, A. Renaux, and G. Ginolhac. Riemannian geometry and Cramér-Rao bound for blind separation of Gaussian sources. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP-2020)*. IEEE, accepted.
- [15] B. Vandereycken, P.-A. Absil, and S. Vandewalle. A Riemannian geometry with complete geodesics for the set of positive semidefinite matrices of fixed rank. *IMA Journal of Numerical Analysis*, 33(2):481–514, 2012.
- [16] R. Bhatia. *Positive definite matrices*. Princeton University Press, 2009.
- [17] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, Princeton, NJ, USA, 2008.
- [18] P.-A. Absil and J. Malick. Projection-like retractions on matrix manifolds. *SIAM Journal on Optimization*, 22(1):135–158, 2012.
- [19] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre. Manopt, a Matlab toolbox for optimization on manifolds. *Journal of Machine Learning Research*, 15:1455–1459, 2014.
- [20] E. Moreau and O. Macchi. New self-adaptive algorithms for source separation based on contrast functions. In *Higher-Order Statistics, 1993., IEEE Signal Processing Workshop on*, pages 215–219. IEEE, 1993.