

---

# Submission and Formatting Instructions for International Conference on Machine Learning (ICML 2016)

---

## Abstract

Independent Component Analysis (ICA) - one of the basic tools in data analysis - aims to find a coordinate system in which the components of the data are independent. In many application the number of sources is unknown and may be less than the number of sensors. In such situation we are looking for so-called non-square mixing matrix. Due to computational constraints, principal component analysis is used for dimension reduction prior to ICA (PCA+ICA), which could remove important information. In this paper we present a two method which are dedicated for determining non-square mixing matrix by fitting non Gaussian densities.

## 1. Introduction

ICA is similar in many aspects to principal component analysis (PCA). In PCA we look for an orthonormal change of basis so that the components are not linearly dependent (uncorrelated). ICA can be described as a search for the optimal basis (coordinate system) in which the components are independent. Let us now, for the readers convenience, describe how the ICA works. The data are represented by the random vector  $x$  and the components as the random vector  $s$ . Our aim is to transform the observed data  $x$  into maximally independent components  $s$  with respect to some measure of independence. Here we use a linear static transformation  $W$ , called the *transformation matrix*, combined with the formula

$$s = Wx.$$

Popular ICA methodology does not directly attempt to find components that are independent but rather components that are as non-Gaussian as possible. This follows from the fact that one of the theoretical foundations of ICA is given by the dual view at the Central Limit Theorem (Hyvärinen & Oja, 2000), which states that the distribution of the sum (average or linear combination) of  $N$  independent random

variables approaches Gaussian as  $N \rightarrow \infty$ . Obviously if all source variables are Gaussian, the ICA method will not work.

There exists many different approaches to ICA which uses negentropy (Hyvärinen & Oja, 2000), cumulant-based methods (Cardoso & Souloumiac, 1993; Virta et al., 2015), maximum likelihood methods (Chen et al., 2006; Samworth et al., 2012) and methods that directly minimize a measure of dependence (Stögbauer et al., 2004; Matteson & Tsay, 2016).

In many application the number of sources is unknown and may be less than the number of sensors. In such situation we are looking for so-called non-square mixing matrix. In practice, PCA is applied to the observations prior to classic ICA (PCA+ICA) to meet the assumption of square mixing and to reduce computational costs (Hyvärinen et al., 2004). PCA+ICA is commonly used to identify brain networks in functional magnetic resonance imaging (fMRI) (Beckmann, 2012; Green et al., 2002) and hyperspectral unmixing (Wang et al., 2015; Caiafa et al., 2008).

The problem in such approach is that interesting independent components (ICs) could be mixed in several principal components that are discarded and then these ICs cannot be recovered.

In the paper we present two methods dedicated to a maximum-likelihood framework. In the first case we are looking directly  $d \leq D$  independent component by maximization of likelihood function. The second method work in full dimensional space by estimating density congaing  $d$  non-gaussian components (independent ones) and  $D - d$  gaussian ones which model a noise.

[!!!Opisac w miare dokladnie nasze podejscie!!!]

## 2. basic tools

### 2.1. Orthogonal projection onto affine subspaces

Suppose that we have an affine subspace generated over  $m \in \mathbb{R}^D$ ,  $V \in \mathbb{R}^{D \times d}$ , where  $V = [v_1, \dots, v_k]$  (or more precisely its consecutive columns) is the base of linear part of  $P$  with  $d$  elements, that is

$$M = m + \text{span}(V) = m + \{Vr : r \in \mathbb{R}^d\} = \{m + r_1 v_1 + \dots + r_d v_d : r_i \in \mathbb{R}\}$$

We are interested in the coordinates of the point  $x \in \mathbb{R}^D$  after the orthogonal projection onto  $P$  with respect to the base. This can restated as the search for coordinates  $r = (r_1, \dots, r_d)^T \in \mathbb{R}^d$  such that

$$r = \operatorname{argmin}_{s \in \mathbb{R}^d} \|x - (m + Vs)\|^2 = \operatorname{argmin}_{s \in \mathbb{R}^d} \|x - (m + s_1 v_1 + \dots + s_d v_d)\|^2.$$

The formula can be obtained by the least squares solution to the problem  $m + Vr = x$ :

$$r_1 v_1 + \dots + r_d v_d = x - m,$$

which is given by:

$$r = (V^T V)^{-1} V^T (x - m) \in \mathbb{R}^d.$$

## 2.2. Integration on subspaces

For the integration over  $C^1$  submanifolds of  $\mathbb{R}^D$  refer the reader to (Munkres, 1997; Federer, 2014). If we are given an a  $C^1$  submanifold  $M$  of dimension  $d$  of  $\mathbb{R}^D$ , then we have a default restriction of Lebesgue measure to  $M$ , which we denote by  $\lambda_d$  (formally, it is the normalization of  $d$ -dimensional Haar measure).

In the case we are interested in, when  $M$  is an affine subspace, to integrate a function over  $M$  we can take a point  $m \in M$  and base  $V$  of the linear part of  $M$ , and then

$$\int_M f(x) d\lambda_d(x) = \det(V^T V)^{1/2} \int_{\mathbb{R}^d} f(m + Vr) d\lambda_d(r).$$

With respect to measure  $\lambda_d$  in  $\mathbb{R}^D$  we can consider the singular densities (that is those defined only on  $M$ , or equivalently zero except for  $M$ ). In the most important case of Gaussian densities, if  $m \in \mathbb{R}^D$  and  $\Sigma$  is a symmetric non-negative matrix with rank  $d$ , then by  $N(m, \Sigma)$  we denote the function with support in  $M = \{m + \Sigma^{1/2} r : r \in \mathbb{R}^D\}$  and the density given by

$$\mathcal{N}(m, \Sigma)(x) = \frac{1}{\sqrt{\det^*(2\pi\Sigma)}} e^{-\frac{1}{2}(x-m)^T \Sigma^\dagger (x-m)} \text{ for } x \in M,$$

where  $\Sigma^\dagger$  is the generalized Moore-Penrose inverse and  $\det^*$  is the pseudo-determinant<sup>1</sup>.

## 2.3. Push-forward of measures

Since we know how to integrate functions on affine subspaces, let us discuss the natural method of defining (by push-forward) measures on such subspaces, for more information see [https://en.wikipedia.org/wiki/Pushforward\\_measure](https://en.wikipedia.org/wiki/Pushforward_measure), <http://www.mat.univie.ac.at/~gerald/>

<sup>1</sup>That is the product of all nonzero eigenvalues.

<ftp://book-fa/index.html> (page 256) and (Bogachev, 2007). We assume as before, that  $M$  is an affine subspace of  $\mathbb{R}^D$  of dimension  $d$ , and that we fix  $m$  (coordinate center) and  $V$  (base of linear part of  $M$ ). Assume that we are given an affine function

$$a : \mathbb{R}^d \ni r \rightarrow m + Vr \in M \subset \mathbb{R}^D.$$

Then  $m$  and  $V$  introduce a coordinate system on  $V$ , with center at  $m$ .

Suppose that we are given a measure  $\mu$  on  $\mathbb{R}^d$  with density  $f$ . Then we can push-forward (transport) the measure  $\mu$  onto  $M$  through the map  $a$  to obtain the measure by the formula

$$(a_*\mu)(B) = \mu(a^{-1}B) \text{ for } B \subset \mathbb{R}^D.$$

By applying the knowledge of integration over submanifolds, we obtain that the measure  $a_*\mu$  with support in  $M$  has the singular density with respect to  $\lambda_d$  given by

$$f_{m,V}(x) = \begin{cases} \frac{1}{\sqrt{\det(V^T V)}} f(a_m^{-1}x) & \text{if } x \in M, \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Roughly speaking the above means that if we have a dataset  $W \subset \mathbb{R}^d$  which comes from the density  $f$  on  $\mathbb{R}^d$ , then  $a_{m,V}(W)$  is clearly supported in  $M$  and comes from the singular density  $f_{m,V}$  given by (1).

In the particular case when  $W$  comes from the normal density  $\mathcal{N}(m_d, \Sigma_d)$  in  $\mathbb{R}^d$ , then  $a_{m,V}(W)$  has the singular normal density  $\mathcal{N}(m + Vm_d, V^T \Sigma_d V)$  in  $\mathbb{R}^D$ .

## 2.4. Measure of nongaussianity

We consider the similar idea to the Kullback-Leibler.

## 2.5. Construction of densities

We can define the family of singular densities on affine subspaces of dimension  $d$ , by taking the transport.

In this subsection we describe the basic construction of product measures and densities. Given functions  $f_1, f_2$  on  $\mathbb{R}^{d_1}, \mathbb{R}^{d_2}$  by

$$(f_1 \otimes f_2)(x_1, x_2) = f_1(x_1) \cdot f_2(x_2) \text{ for } (x_1, x_2) \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$$

we denote the tensor product of  $f_1$  and  $f_2$ . Observe that if  $f_1, f_2$  are densities, then so is  $f_1 \otimes f_2$ .

If  $\mathcal{F}$  is a family of densities on  $\mathbb{R}$ , then by  $\mathcal{F}^{\otimes d}$  ( $d$ -th tensor power of  $\mathcal{F}$ ) we denote the family of densities on  $\mathbb{R}^d$  given by

$$\mathcal{F}^{\otimes k} = \{f_1 \otimes \dots \otimes f_d : f_i \in \mathcal{F}\}.$$

## 2.6. Our case

We assume that we have a family  $\mathcal{F}$  larger than Gaussians on  $\mathbb{R}$ .

We have

$$\text{KL}(X, \text{aff}(\mathcal{S}^{\otimes d}), \mathcal{G}^d) = \inf_{m, V} \text{KL}((v^T V)^{-1} V^T (X - m), \mathcal{S}^{\otimes d}, \mathcal{G}).$$

Notation:  $x[m, V]$ . By the  $i$ -th coordinate we denote  $x[m, V]_i$ .

Thus

$$\text{KL}(X, \text{aff}(\mathcal{S}^{\otimes d}), \mathcal{G}^d) = \inf_{m, V} \left( \sum_{i=1}^d \text{mle}(X[m, V]_i, \mathcal{S}) - \text{mle}(X[m, V], \mathcal{G}) \right), \quad \text{KL}(X, \mathcal{F}, \mathcal{G}) = \text{mle}(X, \mathcal{F}) - \text{mle}(X, \mathcal{G}).$$

where the minus has the direct formula which can be computed.

## 3. First approach: global estimation

We search for the split  $g = f \otimes \mathbb{N}$ , where  $g$  is a normal density on  $\mathbb{R}^{D-d}$ , and  $f$  is  $d$ -dimensional. More precisely, we fix a family  $\mathcal{F}$  of densities on  $\mathbb{R}^d$ , and seek  $m, V$  which maximize the MLE:

$$X \sim a_*(f \otimes g) \text{ for } f \in \mathcal{F}, g \in \mathcal{N}(\mathbb{R}^{D-d}), a = a_{m, V} \in \text{aff}(\mathbb{R}^D)$$

In the case when  $\mathcal{F}$  is one dimensional, the above can be written as:

$$X \sim \det W \cdot f_1(w_1 \circ (x - m)) \cdots f_d(w_d \circ (x - m)) \cdot g_{d+1}(w_{d+1} \circ (x - m))$$

where  $W = [w_1, \dots, w_D] = (V^{-1})^T$ .

## 4. Second approach: projection

We want to find an index which would have the following characteristics:

1. the more non-gaussian data the better,
2. for gaussian data the value zero,
3. invariant under affine transformations.
4.  $k(f * N) < k(f)$  which implies the minimization?

**Theorem 4.1.** ? *Theorem: in the perfect split we obtain original split ?*

*Proof.* We have two random variables which are independent, the second Gaussian. Observe that if the change of coordinates, then sum of independent variables.

We search for minimal entropy (maximal likelihood). Since [Original Entropy Power Inequality]

$$e^{2H(X+Y)} \geq e^{2H(X)} + e^{2H(Y)},$$

and the equality holds only for the gaussians,  $\square$

We propose the possible solution for the ICA. We assume that we are given an affine-invariant family  $\mathcal{F}$  of densities on  $\mathbb{R}^D$ , which contains normal densities  $\mathcal{G}$  (Gaussians). To measure the distance from normality, we define an analogue of Kullback-Leibler divergence [sprawdzic znak, jak entropia to odwrotnie?]:

[czy bierzemy znormalizowane - czy sumaryczne?]

Observe that for a fixed data the second element depends only on the covariance of the data. On the other hand, the first component typically has to be optimized by some gradient methods. Since the formula for the previous part is known

$$\text{mle}(X, \mathcal{G}) = \text{card} X \left( -\frac{1}{2} \ln |\Sigma_X| - \frac{D}{2} \ln(2\pi e) \right).$$

Now consider the situation where we are given a task of finding dimension on possibly smaller space of dimension  $d \leq D$ . In this case assume that we are given a family  $\mathcal{F}^d$  on  $\mathbb{R}^d$ , where  $d \leq D$  (we do not assume that  $\mathcal{F}^d$  is affine invariant, as we obtain it directly from the construction by the fact that we can adapt the base). To fix an affine space  $V$  of dimension  $d$  in  $\mathbb{R}^D$  we choose its center  $m$  and  $d$  linearly independent elements  $V = v_1, \dots, v_d \in \mathbb{R}^D$ .

Now the coordinates<sup>2</sup> in the base  $V$  of orthogonal projection of  $x \in \mathbb{R}^D$  onto  $V$  is given by

$$\lambda_{m, V}^x = (V^T V)^{-1} V^T (x - m) \in \mathbb{R}^d \text{ and } x_{m, v} = m + V \lambda_{m, V}^x. \quad (2)$$

By  $\Lambda_{m, V} = (\lambda_{m, v}^x)$  we denote the coordinates of the whole data set. Now we can project the data to this space, and in those coordinates we can measure the previously defined Kullback-Leibler generalized divergence:

$$(m, V) \rightarrow \text{KL}(\Lambda_{m, V}, \mathcal{F}^d, \mathcal{G}). \quad (3)$$

The minimization of the above function leads to the solution of the ICA problem on the respective subspace.

We will consider it for the family  $\mathcal{F}$  of split Gaussians, however, one can apply any family used in the ICA process.

<sup>2</sup>The formula is the direct consequence of the fact that the orthogonal projection is exactly the solution of least squares solution of the equations  $v\alpha = x - m$ , where  $\alpha = (\alpha_1, \dots, \alpha_d)^T$

It occurs that under weak assumption we can even rank the base vectors of  $V$ . To do so suppose that  $\mathcal{S}^{\otimes d}$  is given as tensor product  $\mathcal{S}^{\otimes d} = \mathcal{S} \otimes \dots \otimes \mathcal{S}$ , where  $\mathcal{S}$  denotes a family of densities on  $\mathbb{R}$  (this is the case of split Gaussians). In other words we assume that every element of  $F \in \mathcal{S}^d$  can be decomposed in the form

$$F(x_1, \dots, x_d) = f_1(x_1) \cdot \dots \cdot f_d(x_d) \text{ where } f_i \in \mathcal{S}.$$

Notation  $\text{aff}(\mathcal{S}^{\otimes d})$  – will denote the space of affine. If we are given a density  $f$  on  $\mathbb{R}^d$ , and an affine map  $A : \mathbb{R}^d \ni \lambda \rightarrow m + V\lambda$ , then the degenerate density on the space  $V$  with respect to the  $d$ -dimensional Lebesgue (Haar) measure  $\lambda_d$  is given by

$$f_V : V \ni x \rightarrow \frac{1}{|A|} f(A^{-1}x)$$

where  $|A|$  is the generalization of determinant given by ... The formula for the KL is therefore given by

$$\sum_x \ln f(A^{-1}p_V x) - \ln N(A^{-1}p_V x).$$

Observe that  $\Sigma \Lambda_V = (A^{-1}p_V) \Sigma (A^{-1}p_V)^T$ . Consequently, the minus part equals

$$\text{card} X \left( -\frac{1}{2} \ln |(A^{-1}p_V) \Sigma (A^{-1}p_V)^T| - \frac{d}{2} \ln(2\pi e) \right).$$

PROCEDURE to compute  $\text{KL}_{m,V}^d(X, \mathcal{S})$ :

- data  $X$  and family of one-dimensional densities on  $\mathcal{S}$  given,
- fix  $m, V$ ,
- put  $\Lambda = (\lambda_{m,V}^x)_{x \in X} \subset \mathbb{R}^d$ ,
- by  $\Lambda_i$  we denote the set consisting of  $i$ -th coordinate of  $\Lambda$ ,
- compute<sup>3</sup>

$$\text{KL}(\Lambda, \mathcal{S}^d, \mathcal{G}) = \sum_{i=1}^d \text{mle}(\Lambda_i, \mathcal{S}) - \text{mle}(\Lambda, \mathcal{G}).$$

We put

$$\text{KL}^d(X, \mathcal{S}) = \inf \text{KL}_{m,V}^d(X, \mathcal{S}).$$

**Theorem 4.2.** *a) Independent of affine transformations b) czy możemy sie zawezic do popdrzestrzeni*

**Problem 4.1.** czy jest znany wzor dla mle przy split gaussian?

**Theorem 4.3.**

<sup>3</sup>sometimes we need optimization

Now suppose that we have found a base  $m, V$  which minimizes (3). Denote by  $(\alpha)_i$   $i$ -th coordinate of  $\alpha$ , then we can rank the vectors according to the non-gaussianity of the  $i$ -th coordinate of the projection:

$$i \rightarrow \text{KL}((X_{m,V})_i, \mathcal{F}^1, \mathcal{G}).$$

We want to introduce a new measure to see if the subspace we found is correct. The model has to be affine independent. To do so, assume that we are given data  $X = (x_i)$  and the transformed/obtained data  $\tilde{X} = (\tilde{x}_i)$ . We define the measure between the best affine transformation between data, to do so by mean squares we solve the problem

$$A\tilde{x}_i + b = x_i.$$

The mean squared error is the desired value:

$$i(X, \tilde{X}) = \frac{1}{N} \sum_{i=1}^N \|x_i - (A\tilde{x}_i + b)\|^2.$$

**Example 4.1.** Take the real-data  $X \subset \mathbb{R}^d$ , add the next  $D - d$  coordinates by some normal density – we obtain new data set  $\tilde{X}$ . Try to find the first  $d$  coordinates.

Measure the value of

$$i(X, \tilde{X}).$$

**Example 4.2.** Take the real-data  $X \subset \mathbb{R}^d$ , add the next  $D - d$  coordinates with zeros. Next perturb all coordinates by some normal density. Try to find the first  $d$  coordinates. Come back by least squares between the original coordinates and the projection.

## 5. przemek

The density of the one-dimensional Split Gaussian distribution is given by the formula

$$SN(x; m, \sigma^2, \tau^2) = \begin{cases} c \cdot \exp[-\frac{1}{2\sigma^2}(x - m)^2], & \text{where } x \leq m \\ c \cdot \exp[-\frac{1}{2\tau^2\sigma^2}(x - m)^2], & \text{where } x > m \end{cases}$$

$$\text{where } c = \sqrt{\frac{2}{\pi}} \sigma^{-1} (1 + \tau)^{-1}.$$

A natural generalization of the univariate split normal distribution to the multivariate settings was presented by (?). Roughly speaking, authors assume that a vector  $x \in \mathbb{R}^d$  follows the multivariate Split Normal distribution, if its principal components are orthogonal and follow the one-dimensional Split Normal distribution.

**Definition 5.1.** A density of the multivariate Split Normal distribution is given by

$$SN_d(x; m, \sigma, \tau) = \prod_{j=1}^d SN(x_j; m_j, \sigma_j^2, \tau_j^2),$$

where  $\mathbf{m} = [m_1, \dots, m_d]^T$ ,  $\sigma = [\sigma_1^2, \dots, \sigma_d^2]^T$  and  $\tau = [\tau_1^2, \dots, \tau_d^2]^T$ .

In our case we will use density on projection on  $d < D$  subspaces. Therefore we need a density  $d$ -subspace Split Normal distribution.

**Definition 5.2.** A density of the multivariate  $d$ -subspace Split Normal distribution is given by

$$SN_{d < D}(\mathbf{x}; \mathbf{m}, W, \sigma^2, \tau^2) = SN_d((W^T W)^{-1} W^T (\mathbf{x} - \mathbf{m}); 0, \sigma^2, \tau^2),$$

where  $(W^T W)^{-1} W^T (\mathbf{x} - \mathbf{m}) \in \mathbb{R}^d$   $\mathbf{w}_j \in \mathbb{R}^D$  is the  $j$ -th column of non-singular matrix  $W = [w_1, \dots, w_d]$ ,  $\mathbf{m} = [m_1, \dots, m_d]^T$ ,  $\sigma = [\sigma_1, \dots, \sigma_d]^T$  and  $\tau = [\tau_1, \dots, \tau_d]^T$ .

Let us recall that the standard Gaussian density in  $\mathbb{R}^d$  is defined by

$$N(\mathbf{x}; \mathbf{m}, \Sigma) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \Sigma^{-1} (\mathbf{x} - \mathbf{m})\right),$$

where  $\mathbf{m}$  denotes the mean,  $\Sigma$  is the covariance matrix.

**Definition 5.3.** A density of the multivariate  $d$ -subspace Normal distribution is given by

$$N_{d < D}(\mathbf{x}; \mathbf{m}, \Sigma, W) = N((W^T W)^{-1} W^T (\mathbf{x} - \mathbf{m}); 0, \Sigma),$$

where  $(W^T W)^{-1} W^T (\mathbf{x} - \mathbf{m}) \in \mathbb{R}^d$   $\mathbf{w}_j \in \mathbb{R}^D$  is the  $j$ -th column of non-singular matrix  $W = [w_1, \dots, w_d]$ ,  $\mathbf{m} = [m_1, \dots, m_d]^T$ ,  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ .

Our goal is to minimize

$$\text{KL}(X, \mathcal{F}, \mathcal{G}) = \text{mle}(X, \mathcal{F}) - \text{mle}(X, \mathcal{G})$$

In our language

$$\begin{aligned} & \text{KL}_{d < D}(X; \mathbf{m}, W, \sigma, \tau, \Sigma) = \\ & = \sum_{\mathbf{x} \in X} \ln(SN_{d < D}(\mathbf{x}; \mathbf{m}, W, \sigma, \tau)) - \sum_{\mathbf{x} \in X} \ln(N_{d < D}(\mathbf{x}; \mathbf{m}, \Sigma, W)) \end{aligned} \quad (4)$$

We known

$$\sum_{\mathbf{x} \in X} \ln(N_{d < D}(\mathbf{x}; \mathbf{m}, \Sigma, W)) = -\frac{d}{2} \ln(2\pi e) - \frac{1}{2} \ln \det(\Sigma_W),$$

where

$$\Sigma_W = \text{cov}(\{(W^T W)^{-1} W^T (\mathbf{x} - \mathbf{m}) : \mathbf{x} \in \mathbb{R}^D\})$$

### 5.1. Optimization problem

The density of the multivariate  $d$ -subspace Normal distribution depends on four parameters  $\mathbf{m} \in \mathbb{R}^d$ ,  $W \in \mathcal{M}(\mathbb{R}^D)$ ,  $\sigma \in \mathbb{R}^d$ ,  $\tau \in \mathbb{R}^d$ . We can find them by minimizing the simpler function, which depends on only  $\mathbf{m} \in \mathbb{R}^d$  and  $W \in \mathcal{M}(\mathbb{R}^D)$ . Other parameters are given by explicit formulas. Let us notice that in this case our minimization problem simplifies to minimizing the function  $\text{mle}(X, \mathcal{F}) = \sum_{\mathbf{x} \in X} \ln(SN_{d < D}(\mathbf{x}; \mathbf{m}, W, \sigma, \tau))$

**Theorem 5.1.** Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be given. Then the likelihood maximized w.r.t.  $\sigma$  and  $\tau$  is

$$\hat{L}(X; \mathbf{m}, W) = \left(\frac{2n}{\pi e}\right)^{dn/2} \left(\prod_{j=1}^d g_j(\mathbf{m}, W)\right)^{-3n/2}, \quad (5)$$

where

$$g_j(\mathbf{m}, W) = s_{1j}^{1/3} + s_{2j}^{1/3},$$

$$s_{1j} = \sum_{i \in I_j} [\mathbf{w}_j^T (\mathbf{x}_i - \mathbf{m})]^2, I_j = \{i = 1, \dots, n : \mathbf{w}_j^T (\mathbf{x}_i - \mathbf{m}) \leq 0\},$$

$$s_{2j} = \sum_{i \in I_j^c} [\mathbf{w}_j^T (\mathbf{x}_i - \mathbf{m})]^2, I_j^c = \{i = 1, \dots, n : \mathbf{w}_j^T (\mathbf{x}_i - \mathbf{m}) > 0\},$$

where  $\omega_j$  is the  $j$ -th column of non-singular matrix  $(W^T W)^{-1} W^T$  and the maximum likelihood estimators of  $\sigma_j^2$  and  $\tau_j$  are

$$\hat{\sigma}_j^2(\mathbf{m}, W) = \frac{1}{n} s_{1j}^{2/3} g_j(\mathbf{m}, W), \quad \hat{\tau}_j(\mathbf{m}, W) = \left(\frac{s_{2j}}{s_{1j}}\right)^{1/3}.$$

*Proof of Theorem 5.1.* Let  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and  $W_\omega = (W^T W)^{-1} W^T$ . We write

$$\mathbf{z}_i = W_\omega (\mathbf{x}_i - \mathbf{m}), \quad \mathbf{z}_{ij} = \omega_j^T (\mathbf{x}_i - \mathbf{m}),$$

for observation  $i$ , where  $i = 1, \dots, n$  and coordinates  $j = 1, \dots, d$ .

Let us consider the likelihood function, i.e.

$$\begin{aligned} L(X; \mathbf{m}, W, \sigma, \tau) &= \prod_{i=1}^n SN_{d < D}(\mathbf{x}_i; \mathbf{m}, W, \sigma, \tau) = \prod_{i=1}^n \prod_{j=1}^d SN(\omega_j^T (\mathbf{x}_i - \mathbf{m}); 0, \sigma_j^2, \tau_j) \\ &= c_1^n \left(\prod_{j=1}^d \sigma_j (1 + \tau_j)\right)^{-n} \prod_{i=1}^n \prod_{j=1}^d \exp\left[-\frac{1}{2\sigma_j^2} z_{ij}^2 (\mathbf{1}_{\{z_{ij} \leq 0\}} + \tau_j \mathbf{1}_{\{z_{ij} > 0\}})\right] \end{aligned}$$

where  $c_1 = \left(\sqrt{\frac{2}{\pi}}\right)^d$ . Now we take the log-likelihood function, i.e.

$$\begin{aligned} & \ln(L(X; \mathbf{m}, W, \sigma, \tau)) \\ &= \ln\left(c_1^n \left(\prod_{j=1}^d \sigma_j (1 + \tau_j)\right)^{-n}\right) + \sum_{i=1}^n \sum_{j=1}^d \left[-\frac{1}{2\sigma_j^2} z_{ij}^2 (\mathbf{1}_{\{z_{ij} \leq 0\}} + \tau_j \mathbf{1}_{\{z_{ij} > 0\}})\right] \\ &= \ln\left(c_1^n \left(\prod_{j=1}^d \sigma_j (1 + \tau_j)\right)^{-n}\right) - \frac{1}{2} \sum_{j=1}^d \left(\sigma_j^{-2} \sum_{i \in I_j} z_{ij}^2 + \frac{\sigma_j^{-2}}{\tau_j} \sum_{i \in I_j^c} z_{ij}^2\right) \\ &= \ln\left(c_1^n \left(\prod_{j=1}^d \sigma_j (1 + \tau_j)\right)^{-n}\right) - \sum_{j=1}^d \frac{1}{2\sigma_j^2} \left(s_{1j} + \frac{1}{\tau_j} s_{2j}\right). \end{aligned}$$

We fix  $\mathbf{m}$ ,  $W$  and maximize the log-likelihood function over  $\tau$  and  $\sigma$ . In such a case we have to solve the following system of equations

$$\frac{\partial \ln(L(X; \mathbf{m}, W, \sigma, \tau))}{\partial \sigma_j} = -\frac{n}{\sigma_j} + \sigma_j^{-3} (s_{1j} + \tau_j^{-2} s_{2j}) = 0,$$

$$\frac{\partial \ln(L(X; \mathbf{m}, W, \sigma, \tau))}{\partial \tau_j} = -\frac{n}{1 + \tau_j} + \frac{s_{2j}}{\tau_j^3 \sigma_j^2} = 0,$$



for  $j = 1, \dots, d$ . By simple calculations we obtain the expressions for the estimators

$$\hat{\sigma}_j^2(\mathbf{m}, W) = \frac{1}{n} s_{1j}^{2/3} g_j(\mathbf{m}, W), \quad \hat{\tau}_j(\mathbf{m}, W) = \left( \frac{s_{2j}}{s_{1j}} \right)^{1/3}.$$

Substituting it into the log-likelihood function, we get

$$\begin{aligned} \hat{L}(\mathbf{m}, W) &= \left( \frac{2}{\pi} \right)^{\frac{dn}{2}} \left( \prod_{j=1}^d \frac{1}{\sqrt{n}} g_j(\mathbf{m}, W)^{\frac{3}{2}} \right)^{-n} e^{-\frac{dn}{2}} \\ &= \left( \frac{2n}{\pi e} \right)^{\frac{dn}{2}} \left( \prod_{j=1}^d g_j(\mathbf{m}, W) \right)^{-\frac{3n}{2}}. \end{aligned}$$

□

Thanks to the above theorem, instead of looking for the maximum of the likelihood function, it is enough to obtain the maximum of the simpler function (5) which depends on two parameters  $\mathbf{m} \in \mathbb{R}^d$  and  $W \in \mathcal{M}(\mathbb{R}^d)$

$$l(X; \mathbf{m}, W) = \prod_{j=1}^d g_j(\mathbf{m}, W) \quad (6)$$

where  $w_j$  stands for the  $j$ -th column of matrix  $W$ . Consequently, maximization of (5) is equivalent to minimization of (6), see the following corollary.

**Corollary 5.1.** *Let  $X \subset \mathbb{R}^d$ ,  $\mathbf{m} \in \mathbb{R}^d$ ,  $W \in \mathcal{M}(\mathbb{R}^d)$  be given, then*

$$\operatorname{argmax}_{\mathbf{m}, W} \hat{L}(X; \mathbf{m}, W) = \operatorname{argmin}_{\mathbf{m}, W} l(X; \mathbf{m}, W).$$

## 5.2. Gradient

One of the possible methods of optimization is the gradient method. Since the minimum of  $l$  is equal to the minimum of  $\ln(l)$ , in this subsection we calculate the gradient of  $\ln(l)$ . Before we prove suitable Theorem 5.2, we recall the following lemma.

**Lemma 5.1.** *Let  $A = (a_{ij})_{1 \leq i, j \leq d}$  be a differentiable map from real numbers to  $d \times d$  matrices then*

$$\frac{\partial \det(A)}{\partial a_{ij}} = \operatorname{adj}^T(A)_{ij}, \quad (7)$$

where  $\operatorname{adj}(A)$  stands for the adjugate of  $A$ , i.e. the transpose of the cofactor matrix.

*Proof.* By the Laplace expansion  $\det A = \sum_{j=1}^d (-1)^{i+j} a_{ij} M_{ij}$  where  $M_{ij}$  is the minor of the entry in the  $i$ -th row and  $j$ -th column. Hence

$$\frac{\partial \det A}{\partial a_{ij}} = (-1)^{i+j} M_{ij} = \operatorname{adj}^T(A)_{ij}.$$

□

Now we are ready to calculate gradient of our cost function.

**Theorem 5.2.** *Let  $X \subset \mathbb{R}^d$ ,  $\mathbf{m} = (\mathbf{m}_1, \dots, \mathbf{m}_d)^T \in \mathbb{R}^d$ ,  $W = (w_{ij})_{1 \leq i, j \leq d}$  non-singular be given. Then  $\nabla_{\mathbf{m}} \ln l(X; \mathbf{m}, W) = \left( \frac{\partial \ln l(X; \mathbf{m}, W)}{\partial \mathbf{m}_1}, \dots, \frac{\partial \ln l(X; \mathbf{m}, W)}{\partial \mathbf{m}_d} \right)^T$ , where*

$$\frac{\partial \ln l(X; \mathbf{m}, W)}{\partial \mathbf{m}_k} = \sum_{j=1}^d \frac{-1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \left( \frac{1}{3s_{1j}^{\frac{2}{3}}} \sum_{i \in I_j} 2w_j^T(x_i - \mathbf{m})w_{jk} + \frac{1}{3s_{2j}^{\frac{2}{3}}} \sum_{i \in I_j^c} 2w_j^T(x_i - \mathbf{m})w_{jk} \right)$$

Moreover,  $\nabla_W \ln l(X; \mathbf{m}, W) = \left[ \frac{\partial \ln \tilde{l}(X; \mathbf{m}, W)}{\partial w_{pk}} \right]_{1 \leq p, k \leq d}$ ,

where

$$\begin{aligned} \frac{\partial \ln \tilde{l}(X; \mathbf{m}, W)}{\partial w_{pk}} &= -\frac{2}{3} (w^{-1})_{pk}^T + \frac{1}{s_{1p}^{\frac{1}{3}} + s_{2p}^{\frac{1}{3}}} \left( \frac{1}{3s_{1p}^{\frac{2}{3}}} \sum_{i \in I_p} 2w_p^T(x_i - \mathbf{m})(x_{ik} - \mathbf{m}_k) \right. \\ &\quad \left. + \frac{1}{3s_{2p}^{\frac{2}{3}}} \sum_{i \in I_p^c} 2w_p^T(x_i - \mathbf{m})(x_{ik} - \mathbf{m}_k) \right). \end{aligned}$$

and

$$\begin{aligned} s_{1j} &= \sum_{i \in I_j} [w_j^T(x_i - \mathbf{m})]^2, \quad I_j = \{i = 1, \dots, n : w_j^T(x_i - \mathbf{m}) \leq 0\}, \\ s_{2j} &= \sum_{i \in I_j^c} [w_j^T(x_i - \mathbf{m})]^2, \quad I_j^c = \{i = 1, \dots, n : w_j^T(x_i - \mathbf{m}) > 0\}. \end{aligned}$$

*Proof of Theorem 5.2.* Let us start with the partial derivative of  $\ln(l)$  with respect to  $\mathbf{m}$ . We have

$$\frac{\partial \ln l(X; \mathbf{m}, W)}{\partial \mathbf{m}_k} = \sum_{j=1}^d \frac{\partial \ln(g_j(\mathbf{m}, W))}{\partial \mathbf{m}_k} = \sum_{j=1}^d \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \frac{\partial (s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}})}{\partial \mathbf{m}_k} = \sum_{j=1}^d \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \left( \frac{1}{3s_{1j}^{\frac{2}{3}}} \sum_{i \in I_j} 2w_j^T(x_i - \mathbf{m})w_{jk} + \frac{1}{3s_{2j}^{\frac{2}{3}}} \sum_{i \in I_j^c} 2w_j^T(x_i - \mathbf{m})w_{jk} \right)$$

Now, we need  $\frac{\partial s_{1j}}{\partial \mathbf{m}_k}$  and  $\frac{\partial s_{2j}}{\partial \mathbf{m}_k}$ , therefore

$$\frac{\partial s_{1j}}{\partial \mathbf{m}_k} = \sum_{i \in I_j} \frac{\partial [w_j^T(x_i - \mathbf{m})]^2}{\partial \mathbf{m}_k} = \sum_{i \in I_j} 2w_j^T(x_i - \mathbf{m}) \frac{\partial w_j^T(x_i - \mathbf{m})}{\partial \mathbf{m}_k} = \sum_{i \in I_j} 2w_j^T(x_i - \mathbf{m})w_{jk}$$

Analogously we get

$$\frac{\partial s_{2j}}{\partial \mathbf{m}_k} = \sum_{i \in I_j^c} -2w_j^T(x_i - \mathbf{m})w_{jk}.$$

Hence

$$\frac{\partial \ln l}{\partial \mathbf{m}_k} = \sum_{j=1}^d \frac{-1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \left( \frac{1}{3s_{1j}^{\frac{2}{3}}} \sum_{i \in I_j} 2w_j^T(x_i - \mathbf{m})w_{jk} + \frac{1}{3s_{2j}^{\frac{2}{3}}} \sum_{i \in I_j^c} 2w_j^T(x_i - \mathbf{m})w_{jk} \right)$$

Now we calculate the partial derivative of  $\ln l(X; \mathbf{m}, W)$  with respect to the matrix  $W$ . We have

$$\frac{\partial \ln l(X; \mathbf{m}, W)}{\partial w_{pk}} = \frac{\partial \ln |\det(W)|^{-\frac{2}{3}}}{\partial w_{pk}} + \sum_{j=1}^d \frac{\partial \ln(g_j(\mathbf{m}, W))}{\partial w_{pk}}.$$

To calculate the derivative of the determinant we use Jacobi's formula (see Lemma 5.1). Hence

$$\begin{aligned} \frac{\partial \ln(\det(W)^{-\frac{2}{3}})}{\partial w_{pk}} &= \det(W)^{\frac{2}{3}} \left( -\frac{2}{3} \right) \det(W)^{-\frac{5}{3}} \frac{\partial \det(W)}{\partial w_{pk}} = -\frac{2}{3} \det(W)^{-1} \frac{\partial \det(W)}{\partial w_{pk}} \\ &= -\frac{2}{3} \frac{1}{\det(W)} \left[ \det(W) (W^{-1})_{pk}^T \right] = -\frac{2}{3} (w^{-1})_{pk}^T, \end{aligned}$$

□

where  $(w^{-1})_{pk}^T$  is the element in the  $p$ -th row and  $k$ -th column of the matrix  $(W^{-1})^T$ . Now we calculate

$$\frac{\partial \ln(g_j(m, W))}{\partial w_{pk}} = \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \frac{\partial (s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}})}{\partial w_{pk}} = \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \left( \frac{1}{3s_{1j}^{\frac{2}{3}}} \frac{\partial s_{1j}}{\partial w_{pk}} + \frac{1}{3s_{2j}^{\frac{2}{3}}} \frac{\partial s_{2j}}{\partial w_{pk}} \right)$$

where

$$\frac{\partial s_{1j}}{\partial w_{pk}} = \sum_{i \in I_j} \frac{\partial [w_j^T(x_i - m)]^2}{\partial w_{pk}} = \sum_{i \in I_j} 2w_j^T(x_i - m) \frac{\partial w_j^T(x_i - m)}{\partial w_{pk}} = \begin{cases} 0, & \text{if } j \neq p \\ \sum_{i \in I_p} 2w_p^T(x_i - m)(x_{ik} - m_k), & \text{if } j = p \end{cases}$$

and  $x_{ik}$  is the  $k$ -th element of the vector  $x_i$ . Analogously we get

$$\frac{\partial s_{2j}}{\partial w_{pk}} = \begin{cases} 0, & \text{if } j \neq p \\ \sum_{i \in I_p} 2w_p^T(x_i - m)(x_{ik} - m_k), & \text{if } j = p. \end{cases}$$

Hence we obtain

$$\frac{\partial \ln l}{\partial w_{pk}} = -\frac{2}{3}(w^{-1})_{pk}^T + \frac{1}{s_{1p}^{\frac{1}{3}} + s_{2p}^{\frac{1}{3}}} \left( \frac{1}{3}s_{1p}^{-\frac{2}{3}} \sum_{i \in I_p} 2w_p^T(x_i - m)(x_{ik} - m_k) + \frac{1}{3}s_{2p}^{-\frac{2}{3}} \sum_{i \in I_p} 2w_p^T(x_i - m)(x_{ik} - m_k) \right).$$

□

## 6. MODEL II

**Definition 6.1.** A density of the multivariate Split Normal  $d$  and Normal  $D - d$  distribution is given by

$$SN_d N_{D-d}(x; m, W, \sigma^2, \tau^2) = \det(W) \prod_{j=1}^d SN(w_j^T(x - m); 0, \sigma_j^2) \prod_{j=d+1}^D N(w_j^T(x - m); 0, \tau_j^2)$$

where  $w_j$  is the  $j$ -th column of non-singular matrix  $W$ ,  $m = (m_1, \dots, m_d)^T$ ,  $\sigma = (\sigma_1, \dots, \sigma_d)$  and  $\tau = (\tau_1, \dots, \tau_{D-d})$ .

## References

- Beckmann, Christian F. Modelling with independent components. *Neuroimage*, 62(2):891–901, 2012.
- Bogachev, Vladimir I. *Measure theory*, volume 1. Springer Science & Business Media, 2007.
- Caiafa, Cesar F, Salerno, Emanuele, Proto, Araceli N, and Fiumi, L. Blind spectral unmixing by local maximization of non-gaussianity. *Signal Processing*, 88(1):50–68, 2008.
- Cardoso, Jean-François and Soudoumiac, Antoine. Blind beamforming for non-gaussian signals. In *IEEE Proceedings F (Radar and Signal Processing)*, volume 140, pp. 362–370. IET, 1993.

- Chen, Aiyu, Bickel, Peter J, et al. Efficient independent component analysis. *The Annals of Statistics*, 34(6):2825–2855, 2006.

- Federer, Herbert. *Geometric measure theory*. Springer, 2014.

- Green, Christopher G, Nandy, Rajesh R, and Cordes, Dietmar. Pca-preprocessing of fmri data adversely affects the results of ica. In *Proceedings of international society of magnetic resonance in medicine*, volume 10, 2002.

- Hyvärinen, Aapo and Oja, Erkki. Independent component analysis: algorithms and applications. *Neural networks*, 13(4):411–430, 2000.

- Hyvärinen, Aapo, Karhunen, Juha, and Oja, Erkki. *Independent component analysis*, volume 46. John Wiley & Sons, 2004.

- Langley, P. Crafting papers on machine learning. In Langley, Pat (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.

- Matteson, David S and Tsay, Ruey S. Independent component analysis via distance covariance. *Journal of the American Statistical Association*, (just-accepted):1–38, 2016.

- Munkres, James R. *Analysis on manifolds*. Westview Press, 1997.

- Samworth, Richard J, Yuan, Ming, et al. Independent component analysis via nonparametric maximum likelihood estimation. *The Annals of Statistics*, 40(6):2973–3002, 2012.

- Stögbauer, Harald, Kraskov, Alexander, Astakhov, Sergey A, and Grassberger, Peter. Least-dependent-component analysis based on mutual information. *Physical Review E*, 70(6):066123, 2004.

- Virta, Joni, Nordhausen, Klaus, and Oja, Hannu. Joint use of third and fourth cumulants in independent component analysis. *arXiv preprint arXiv:1505.02613*, 2015.

- Wang, Nan, Du, Bo, Zhang, Liangpei, and Zhang, Lifu. An abundance characteristic-based independent component analysis for hyperspectral unmixing. *IEEE Transactions on Geoscience and Remote Sensing*, 53(1):416–428, 2015.