

ICA for subspaces

P. Spurek

J. Tabor

P. Rola

Abstract

In its basic form Independent Component Analysis (ICA) aims to find an invertible linear transformation of the data so that the resulting data has independent components. In the case when number of sources is less than the number of sensors, the task has an additional step which consists of filtering out the possible noise components. In such situation we are looking for so-called non-square mixing matrix.

Due to computational constraints, principal component analysis is often used for dimension reduction prior to ICA (PCA+ICA). However, such approach often removes important information. In this paper we present a density based method $\text{ICA}_{\mathcal{SV}}$ based on split-gaussians which is dedicated for determining non-square mixing matrix in ICA maximum likelihood framework.

1 Introduction

Independent component analysis (ICA) is similar in many aspects to principal component analysis (PCA). In PCA we look for an orthonormal base in which the data components are not linearly dependent (uncorrelated), while in ICA we search for the coordinate system in which the components are independent. More precisely the aim of ICA is to transform the observed data \mathbf{X} into maximally independent components \mathbf{S} with use of an invertible linear transformation W , called the *transformation matrix*:

$$\mathbf{S} = W^T \mathbf{X}.$$

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 [16], 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.

Another common approach to ICA based on the maximum likelihood estimation [24] is recently gaining popularity [15, 27, 29]. Then we search for the optimally fitted to data coordinate system B and marginal densities f_i such that the data density factors in base B as the product of marginal densities. To obtain an efficient method and avoid overfitting we have to restrict the marginal densities f_i to a class \mathcal{F} of densities which has not too many parameters which can be easily estimated (clearly from obvious reasons this class has to be different from gaussians). As \mathcal{F} we typically choose the super-Gaussian logistic density or other heavy tails distributions.

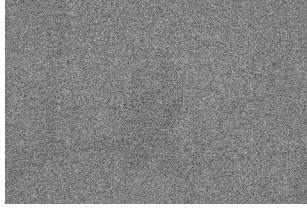
In many applications of ICA we deal with the case when several sensors measure the latent variables and the rest of them record only the noise. This happens when the number of sources is unknown and may be less than the number of sensors (then we are looking for so-called *non-square mixing matrix* W). Such a case is common for example in the identification of brain networks in functional magnetic resonance imaging (fMRI) [3, 10]. In practice, most approaches deal with this problem by first applying PCA to the observations prior to classic ICA (PCA+ICA) to meet the assumption of square mixing and to reduce computational costs [15]. Although numerically effective,



(a) Original images 42049 and 220075.



(b) $ICA_{\mathcal{SN}}$.



(c) FastICA.



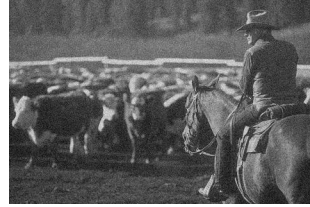
(d) ProDenICA.



(e) $ICA_{\mathcal{SN}}$.



(f) FastICA.



(g) ProDenICA.

Figure 1: Comparison of images separation by our method ($ICA_{\mathcal{SN}}$), with FastICA and ProDenICA. Before mixing by a linear matrix, we added to the first two components given in (a) the third component given by random normal noise. As we see $ICA_{\mathcal{SN}}$ was able to perfectly recover the two first components.

this approach may fail as it is not invariant with respect to linear transformation, since PCA will find a “noise” component if it is sufficiently large.

The aim of this paper is to propose a new density based approach to deal with this case which does not have the above mentioned disadvantage. Our idea is to join the two earlier mentioned approaches to solving ICA - one based on the search for non-gaussian components and the other based on density estimation - to deal with the case when the number of sources is smaller than that of sensors. Observe that the noise typically occurs as a sum of many independent factors, and consequently thanks to the central limit theorem it typically has approximately gaussian distribution. This is why one typically makes the assumption [9] that

noise components are coming from a gaussian noise.

Following the density approach to filter them out we fit the first d -components from a class of \mathcal{F} of densities which is broader than gaussians, while the rest from the gaussians \mathcal{N} (the final choice of the value of parameter $d \in \{1, \dots, D\}$ can be decided by applying either AIC or BIC criterion). Following [29] as \mathcal{F} we take the class \mathcal{SN} of split-normal densities.

Our experiments show, which is illustrated by Figure 1, that $ICA_{\mathcal{SN}}$ works as desired and effectively removes the components which contains gaussian noise. However, we cannot objectively conclude that it is better as compared to other state-of-the-art approaches, since the experiment was conducted in the setting optimal to our method as we assumed that the noise was gaussian.

2 Related works

In literature exist a few approaches dedicated for non-square ICA problem. Most of such method are dedicated for individual methods. Attias and Schreiner [2] derived a likelihood based algorithm for separation of general sequences with a frequency domain implementation. Belouchrani and Cardoso [6] presented a general likelihood approach allowing for additive noise and for non-square mixing matrices. They applied the method to separation of sources taking discrete values and estimated the mixing matrix using an Expectation-Maximization (EM) approach with both a deterministic and a stochastic formulation. In [23] authors used the EM approach for separation of autocorrelated sequences in presence of noise and explored a family of flexible source signal priors based on Gaussian Mixtures.

The assumption is that of square mixing is mostly unrealistic in the case of EEG and FMRI where the number of sources is less than the number of electrodes [4, 26, 28]. Therefore, many of algorithms dedicated for this task use a probabilistic ICA [30]. The noisy ICA model can be approximated using a variant of PCA+ICA [4], where probabilistic PCA is used to estimate the number of components and achieve dimension reduction [30]. In [1] authors developed stochastic EM algorithms to estimate the noisy model and proposed parametric methods.

Other methods exploring non-Gaussian structure in multivariate data include non-Gaussian component analysis (NGCA) and projection pursuit [7, 19]. NGCA is a more general case of linear non-Gaussian component analysis (LNGCA) [25] that allows non-linear dependence between the non-Gaussian components.

In the paper [22] authors propose a novel adaptive twostage deflation-based FastICA algorithm that allows one to use different nonlinearities for different components and optimizes the order in which the components are extracted.

3 Theoretical foundations of ICA_{SN}

In this section we present theoretical foundations of the method. We begin with the statement of the problem, next we focus our attention on the presentation of the class of densities we discuss. Last we show the gradient of the method which is needed in the optimization procedure.

3.1 Statement of the problem

Since this general idea of the search for ICA with the use of maximum likelihood is essential in our further considerations, for the convenience of the reader we first describe it briefly. Assume that the random vector \mathbf{X} in \mathbb{R}^D has the density function $F(\mathbf{x})$. Suppose that the components of \mathbf{X} are not independent, but that we know (or suspect) that there is a basis B (we put $W^T = B^{-1}$) such that in that base the components of \mathbf{X} become independent. Observe that where $\omega_i^T \mathbf{x}$ is the i -th coefficient of \mathbf{x} in the basis B (ω_i denotes the i -th column of W), and therefore there exist densities f_1, \dots, f_D such that

$$F(\mathbf{x}) = \det(W) \cdot f_1(\omega_1^T \mathbf{x}) \cdot \dots \cdot f_d(\omega_d^T \mathbf{x}). \quad (1)$$

Given W and densities $(f_i)_{i=1}^D$ we introduce notation to represent RHS of the above equation:

$$F_W(f_1, \dots, f_D)(\mathbf{x}) = \det(W) \cdot f_1(\omega_1^T \mathbf{x}) \cdot \dots \cdot f_d(\omega_d^T \mathbf{x}).$$

Thus we may state the density based formulation of ICA in the case we have only a sample X from random vector \mathbf{X} .

GENERAL ICA PROBLEM (maximum likelihood formulation).

Find densities f_i and matrix W , so that F given by (1) optimally fits the data $X = (\mathbf{x}_i)$ with respect

to the likelihood, that is that the value

$$\sum_i \log F_W(f_1, \dots, f_D)(\mathbf{x}_i)$$

is maximized.

Since the search over the space of all densities is not feasible, and could lead to overfitting, we naturally have to reduce to a subclass of all densities on \mathbb{R} parametrized by a finite amount of parameters. Clearly, since ICA does not work if the data are gaussian, we have to choose a family \mathcal{F} of densities which is distant from Gaussian ones.

ICA FOR DENSITY CLASS \mathcal{F} .

Find matrix W and densities

$$f_1, \dots, f_D \in \mathcal{F},$$

such that the value of

$$\sum_i \log F_W(f_1, \dots, f_D)(\mathbf{x}_i)$$

is maximized.

Similarly to [29] as a class \mathcal{F} we are going to take the class of split-gaussians, as as it is easy to deal with (small number of parameters) and is resistant to outliers¹.

As mentioned in the introduction, we assume that components which we would like to filter-out, are coming from a gaussian noise, and the aim it to fit the first d -components from a larger class of densities, while the rest from the gaussians \mathcal{N} . Thus our final problem can be stated as follows.

ICA FOR DENSITY CLASS \mathcal{F} WITH d SOURCES.

Find matrix W , densities

$$f_1, \dots, f_d \in \mathcal{F} \text{ and normal densities } f_{d+1}, \dots, f_D \in \mathcal{N},$$

so that the value of

$$\sum_i \ln F_W(f_1, \dots, f_D)(\mathbf{x}_i)$$

is maximized.

Observe that the solution to the above problem is linearly invariant, that is if W is optimal for X an A is linear, then W_A is optimal for AX , where $W_A = (A^{-1})^T W$.

The continuous version of the condition we maximize in the case we know the density f of the random variable \mathbf{X} limits to

$$\int \ln F_W(f_1, \dots, f_D)(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = -H(f, F_W(f_1, \dots, f_D)),$$

where the cross entropy $H(f, g)$ is given by the sum of entropy $H(f)$ and Kullback-Leibler divergence $D_{KL}(f, g)$. Thus the continuous version of the ICA problem with d sources reduces to the minimization of

$$D_{KL}(f, F_W(f_1, \dots, f_D))$$

over all matrices W and densities $f_1, \dots, f_D \in \mathcal{F}$. Since for fixed f Kullback-Leibler divergence is minimized for $g = f$, we arrive at the following result, which says that in the ideal case by the discussed approach we restore the unmixing matrix if it exists.

¹The reason is that split gaussians, instead at fitting the distribution with respect to heavy tails, fits the asymmetry of the data.

Theorem 3.1. *Let F be a density such that there exist matrix \overline{W} and densities*

$$\hat{f}_1, \dots, \hat{f}_d \in \mathcal{F} \text{ and } \hat{f}_{d+1}, \dots, \hat{f}_D \in \mathcal{N}$$

such that

$$F = F_{\overline{W}}(\hat{f}_1, \dots, \hat{f}_D).$$

Then

$$\overline{W}, \hat{f}_1, \dots, \hat{f}_D = \operatorname{argmin}\{F_W(f_1, \dots, f_D) : W, f_1, \dots, \bar{f}_d \in \mathcal{F}, f_{d+1}, \dots, f_D \in \mathcal{N}\}.$$

3.2 Split normal distribution

In this section we discuss the class \mathcal{F} we will use in our final algorithm of $\text{ICA}_{\mathcal{SN}}$. The density of \mathcal{SN} , the one-dimensional split normal distribution [31], is given by the formula

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

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

As we see the split normal distribution comes from merging two opposite halves of two normal distributions in their common mode. The main advantage of split normal distributions over normal one is that it allows data asymmetry. In 1982 John [17] showed that the likelihood function can be expressed in a form in which the scale parameters σ and τ are an explicit function of the location parameter m . In the case when $\mathcal{F} = \mathcal{SN}$ the density class considered in the previous subsection is given in the explicit form by the following observation.

Observation 3.1. *A density of the multivariate split normal d and normal $D - d$ distribution is given by*

$$\mathcal{SN}_d \mathcal{N}_{D-d}(\mathbf{x}; \mathbf{m}, W, \sigma^2, \tau^2) = \det(W) \prod_{j=1}^d \mathcal{SN}(\omega_j^T(\mathbf{x} - \mathbf{m}); 0, \sigma_j^2, \tau_j^2) \cdot \prod_{j=d+1}^D \mathcal{N}(\omega_j^T(\mathbf{x} - \mathbf{m}); 0, \sigma_j^2),$$

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

Observe that the above density probability function has mode in \mathbf{m} . As a consequence of result of John [17] we can maximize the likelihood of the above function on data X with respect to σ and τ .

Theorem 3.2. *Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be given, and let $\mathbf{m} \in \mathbb{R}^D$ and matrix W be fixed. Then the likelihood maximized w.r.t. σ and τ is*

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

where

$$\begin{aligned} g_j(\mathbf{m}, W) &= s_{1j}^{1/3} + s_{2j}^{1/3}, \\ s_{1j} &= \sum_{i \in I_j} [\omega_j^T(\mathbf{x}_i - \mathbf{m})]^2, I_j = \{i: \omega_j^T(\mathbf{x}_i - \mathbf{m}) \leq 0\}, \\ s_{2j} &= \sum_{i \in I_j^c} [\omega_j^T(\mathbf{x}_i - \mathbf{m})]^2, I_j^c = \{i: \omega_j^T(\mathbf{x}_i - \mathbf{m}) > 0\}, \end{aligned}$$

and the maximum likelihood estimators of σ_j^2 and τ_j are

$$\begin{aligned}\hat{\tau}_j(\mathbf{m}, W) &= \left(\frac{s_{2j}}{s_{1j}} \right)^{1/3}, \quad 1 \leq j \leq d \\ \hat{\sigma}_j^2(\mathbf{m}, W) &= \begin{cases} \frac{1}{n} s_{1j}^{2/3} g_j(\mathbf{m}, W), & 1 \leq j \leq d \\ \frac{1}{n} (s_{1j} + s_{2j}), & d < j \leq D. \end{cases}\end{aligned}\quad (3)$$

Proof. See Section 5 (Appendix A). \square

Thanks to the above theorem we can reduce the search for the maximum of the log-likelihood function for two parameters $\mathbf{m} \in \mathbb{R}^d$ and $W \in \mathcal{M}(\mathbb{R}^d)$.

$$l(X; \mathbf{m}, W) = \frac{1}{|\det(W)|^{2/3}} \prod_{j=1}^d g_j(\mathbf{m}, W) \prod_{j=d+1}^D (s_{1j} + s_{2j})^{1/3} \quad (4)$$

where w_j stands for the j -th column of matrix W . Consequently, maximization of likelihood function is equivalent to minimization of $\ln l$.

Corollary 3.1. *Let $X \subset \mathbb{R}^d$, $\mathbf{m} \in \mathbb{R}^d$, $W \in \mathcal{M}(\mathbb{R}^d)$ be given, then*

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

To minimize $\ln l$ with the use classical gradient descent method, we need the formula for $\nabla \ln l$ (gradient of the cost function).

Theorem 3.3. *Let $X \subset \mathbb{R}^d$, $\mathbf{m} = (m_1, \dots, m_d)^T \in \mathbb{R}^d$, $W = (\omega_{ij})_{1 \leq i, j \leq d}$ non-singular be given. Then $\nabla_{\mathbf{m}} \ln l(X; \mathbf{m}, W) = \left(\frac{\partial \ln l(X; \mathbf{m}, W)}{\partial m_1}, \dots, \frac{\partial \ln l(X; \mathbf{m}, W)}{\partial m_d} \right)^T$, where*

$$\begin{aligned}\frac{\partial \ln l(X; \mathbf{m}, W)}{\partial m_k} &= \sum_{j=1}^d \frac{-2}{3(s_{1j}^{1/3} + s_{2j}^{1/3})} \left(\frac{1}{s_{1j}^{2/3}} \sum_{i \in I_j} \omega_j^T (\mathbf{x}_i - \mathbf{m}) \omega_{jk} + \right. \\ &\quad \left. \frac{1}{s_{2j}^{2/3}} \sum_{i \in I_j^c} \omega_j^T (\mathbf{x}_i - \mathbf{m}) \omega_{jk} \right) + \sum_{j=d+1}^D \frac{-2}{3(s_{1j} + s_{2j})} \cdot \\ &\quad \left(\sum_{i \in I_j} \omega_j^T (\mathbf{x}_i - \mathbf{m}) \omega_{jk} + \sum_{i \in I_j^c} \omega_j^T (\mathbf{x}_i - \mathbf{m}) \omega_{jk} \right).\end{aligned}$$

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

$$\begin{aligned}& -\frac{2}{3} (\omega^{-1})_{pk}^T + \frac{2}{3(s_{1p}^{1/3} + s_{2p}^{1/3})} \left(s_{1p}^{-2/3} \sum_{i \in I_p} \omega_p^T (\mathbf{x}_i - \mathbf{m}) (\mathbf{x}_{ik} - \mathbf{m}_k) \right. \\ & \quad \left. + s_{2p}^{-2/3} \sum_{i \in I_p^c} \omega_p^T (\mathbf{x}_i - \mathbf{m}) (\mathbf{x}_{ik} - \mathbf{m}_k) \right) + \frac{2}{3(s_{1p} + s_{2p})} \left(\sum_{i \in I_p} \omega_p^T (\mathbf{x}_i - \mathbf{m}) (\mathbf{x}_{ik} - \mathbf{m}_k) \right. \\ & \quad \left. + \sum_{i \in I_p^c} \omega_p^T (\mathbf{x}_i - \mathbf{m}) (\mathbf{x}_{ik} - \mathbf{m}_k) \right)\end{aligned}$$

and

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

Proof. See Section 6 (Appendix B). \square

Thanks to the above Theorem we are able to use in our experiments the gradient descent for finding the minimum of our cost function.

Table 1: Tucker’s congruence coefficients for reconstruction of two images.

	ICA _{SV}	FastICA LOGCOSH	FastICA EXP	FastICA KURTOSIS	INFOMAX TANH	INFOMAX TANGENT	INFOMAX LOGISTIC	JADE	PEARSONICA	ProDENICA	FixNA
4.1.01	0.99984	0.59585	0.59967	0.58722	0.59596	0.58965	0.59589	0.58333	0.46302	0.83162	0.99948
4.1.02	1	0.93572	0.93315	0.94044	0.93564	0.93925	0.93569	0.94215	0.92294	0.98521	0.99325
4.1.06	0.9959	0.61189	0.61636	0.59603	0.61599	0.60059	0.69578	0.5629	0.5971	0.99852	0.99937
4.1.03	0.8676	0.79174	0.78969	0.7982	0.78986	0.79647	0.72715	0.80803	0.79781	0.87541	0.99811
4.2.04	0.99955	0.60087	0.6007	0.60071	0.60088	0.60094	0.60087	0.59933	0.57862	0.98594	0.88695
5.2.10	0.99974	0.94413	0.9406	0.95526	0.94426	0.94951	0.94395	0.96464	0.98256	0.97757	0.88972
4.2.02	0.99679	0.59655	0.59712	0.59554	0.59658	0.59538	0.59656	0.59434	0.59743	0.99237	0.96088
5.2.08	0.99021	0.97131	0.96951	0.97395	0.97122	0.97433	0.97128	0.97647	0.96843	0.98366	0.94757
BOAT.512	0.99732	0.58702	0.58689	0.58784	0.58692	0.58748	0.58623	0.5865	0.58728	0.99835	0.93781
5.3.01	0.80906	0.98201	0.98196	0.98213	0.98197	0.98213	0.98163	0.98177	0.98209	0.97381	0.97218
ELAINE.512	0.96643	0.58926	0.58919	0.59017	0.58922	0.58987	0.58935	0.58929	0.58901	0.65476	0.71303
4.2.03	0.99839	0.96629	0.96628	0.96641	0.96628	0.96638	0.96631	0.9663	0.96624	0.99326	0.79288
119082	0.99995	0.59432	0.59502	0.59347	0.59391	0.59374	0.59492	0.58197	0.58604	0.99233	0.76561
157055	0.99804	0.94331	0.94273	0.94397	0.94363	0.94377	0.94281	0.94928	0.94804	0.93969	0.73722
42049	0.9999	0.63238	0.61046	0.62306	0.63239	0.62743	0.6324	0.61757	0.5872	0.93106	0.88476
220075	0.99845	0.79121	0.70406	0.91538	0.79133	0.90155	0.79148	0.92798	0.96143	0.9678	0.93793
43074	0.98658	0.58577	0.58375	0.59414	0.58555	0.59142	0.58612	0.59278	0.55845	0.57092	0.78144
295087	0.99757	0.97911	0.98011	0.97119	0.97923	0.97465	0.97891	0.97306	0.97978	0.99416	0.80112
38092	0.95367	0.58677	0.58678	0.58699	0.58677	0.58692	0.58675	0.587	0.58572	0.38511	0.76369
167062	0.99999	0.9916	0.9916	0.99139	0.9916	0.99148	0.99161	0.99137	0.99148	0.99881	0.74187

4 Experiments

To compare ICA_{SV} to other state-of-the-art approaches we use Tucker’s congruence coefficient [20] which values range between -1 and $+1$. It can be used to study the similarity of extracted factors across different samples. Generally, a congruence coefficient of 0.9 indicates a high degree of factor similarity, while a coefficient of 0.95 or higher indicates that the factors are virtually identical.

We evaluate our method in the context of 2D and hyperspectral images. For comparison we use R package `ica` [12], `PearsonICA` [18], `ProDenICA` [11], `tsBSS` [21]. The most popular method used in practice is FastICA [14, 13] algorithm, which uses negentropy. In this context we can use three different functions to estimate neg-entropy: logcosh, exp and kurtosis. We also compare our method with algorithm using Information-Maximization (Infomax) approach [5]. Similarly to FastICA we consider three possible non-linear functions: hyperbolic tangent, logistic and extended Infomax.

4.1 Separation of images

One of the most popular application of ICA is the separation of images. In our experiments we use four images from the USC-SIPI Image Database of size 256×256 pixels (4.1.01, 4.1.06, 4.1.02, 4.1.03) and eight of size 512×512 pixels (4.2.04, 4.2.02, boat.512, elaine.512, 5.2.10, 5.2.08, 5.3.01, 4.2.03). We also use 8 images from the Berkeley Segmentation Dataset of size 482×321 with indexes (#119082, #42049, #43074, #38092, #157055, #220075, #295087, #167062).

We make random pairs of above images and one component with noise (random sample from Gaussian distribution $\mathcal{N}(0,1)$) and use them as a source signal combined by the mixing matrix

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & -1 \\ -1 & 1 & -1 \end{bmatrix}. \text{ Our goal was to reconstruct two original images by using only the knowledge}$$

about mixed ones. The visualization of this process we present in Fig. 1. The results of this experiment are presented in Tab. 1 where we present Tucker’s congruence coefficients which shows that almost in all cases ICA_{SV} obtains best results. This is illustrated in Figure 1, where we can see that ICA_{SV} almost perfectly recovered source signal. Although this is not surprising as the experiments were in fact conducted in the setting which favored our approach, as we chose the



(a) Ground truth layers which contains 4 channels: #1 Asphalt, #2 Grass, #3 Tree and #4 Roof.



(b) The effect of the ICA_{SN} method.



(c) The effect of the FastICA (logcosh) method.



(d) The effect of the ProDenICA method.

Figure 2: Results of image separation with the uses of various ICA algorithms.

noise to be gaussian, this shows that ICA_{SN} works as desired and deals well with removing gaussian components from the data.

4.2 Hyperspectral Unmixing

Independent component analysis has been recently applied into hyperspectral unmixing [32, 8] as a result of its low computation time and its ability to perform without prior information. In this subsection we apply simple example which suggests that our method also can by used for spectral data.

Urban data [33, 35, 34] is one of the most widely used hyperspectral data-sets used in the hyperspectral unmixing study. Each image has 307×307 pixels, each of which corresponds to a 2×2 m area. In this image, there are 210 wavelengths ranging from 400 nm to 2500 nm, resulting in a spectral resolution of 10 nm. After the channels 1–4, 76, 87, 101–111, 136–153 and 198–210 are

removed (due to dense water vapor and atmospheric effects), there remain 162 channels (this is a common preprocess for hyperspectral unmixing analyses). There is ground truth [33, 35, 34], which contains 4 channels: #1 Asphalt, #2 Grass, #3 Tree and #4 Roof.

A highly mixed area is cut from the original data set in this experiment (similar example was showed in [32]), with the size of 200×150 pixels.

In our experiment we compared $\text{ICA}_{\mathcal{SV}}$ to other two popular ICA methods – ProDenICA and FastICA, see Fig. 2. Observe that $\text{ICA}_{\mathcal{SV}}$ and ProDenICA give layers which seem to contain more information than FastICA, as the last component in FastICA contains mainly noise.

5 Appendix A

Proof of Theorem 3.2. Let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. We write

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

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

Let us consider the likelihood function, i.e.

$$\begin{aligned} L(X; \mathbf{m}, W, \sigma, \tau) &= \prod_{i=1}^n S N_d N_{D-d}(\mathbf{x}_i; \mathbf{m}, W, \sigma^2, \tau^2) \\ &= \prod_{i=1}^n |\det(W)| \prod_{j=1}^d S N(\omega_j^T(\mathbf{x}_i - \mathbf{m}); 0, \sigma_j^2, \tau_j^2) \cdot \\ &\quad \prod_{j=d+1}^D N(\omega_j^T(\mathbf{x}_i - \mathbf{m}); 0, \sigma_j^2) = \left(c_1 |\det(W)|\right)^n \\ &\quad \left(\prod_{j=1}^d \sigma_j (1 + \tau_j)\right)^{-n} \prod_{i=1}^n \prod_{j=1}^d \exp\left[-\frac{1}{2\sigma_j^2} z_{ij}^2 (\mathbb{1}_{\{z_{ij} \leq 0\}} + \right. \\ &\quad \left. \tau_j^{-2} \mathbb{1}_{\{z_{ij} > 0\}})\right] \left(\prod_{j=d+1}^D \sigma_j\right)^{-n} \prod_{i=1}^n \prod_{j=d+1}^D \exp\left[-\frac{1}{2\sigma_j^2} z_{ij}^2\right], \end{aligned}$$

where $c_1 = \left(\sqrt{\frac{2}{\pi}}\right)^d \cdot \left(\frac{1}{\sqrt{2\pi}}\right)^{D-d}$. Now we take the log-likelihood function, i.e. $\ln(L(X; \mathbf{m}, W, \sigma, \tau)) =$

$$\begin{aligned} &\ln \left(\left(c_1 |\det(W)|\right)^n \left(\prod_{j=1}^d \sigma_j (1 + \tau_j)\right)^{-n} \left(\prod_{j=d+1}^D \sigma_j\right)^{-n} \right) + \\ &\sum_{i=1}^n \sum_{j=1}^d \left[-\frac{1}{2\sigma_j^2} z_{ij}^2 (\mathbb{1}_{\{z_{ij} \leq 0\}} + \tau_j^{-2} \mathbb{1}_{\{z_{ij} > 0\}}) \right] + \\ &\sum_{i=1}^n \sum_{j=d+1}^D \left[-\frac{1}{2\sigma_j^2} z_{ij}^2 \right] \\ &= \ln \left(\left(c_1 |\det(W)|\right)^n \left(\prod_{j=1}^d (1 + \tau_j)\right)^{-n} \left(\prod_{j=1}^D \sigma_j\right)^{-n} \right) - \\ &\frac{1}{2} \sum_{j=1}^d \left(\sigma_j^{-2} \sum_{i \in I_j} z_{ij}^2 + \frac{\sigma_j^{-2}}{\tau_j^2} \sum_{i \in I_j^c} z_{ij}^2 \right) - \\ &\frac{1}{2} \sum_{j=d+1}^D \sigma_j^{-2} \left(\sum_{i \in I_j} z_{ij}^2 + \sum_{i \in I_j^c} z_{ij}^2 \right) \\ &= \ln \left(\left(c_1 |\det(W)|\right)^n \left(\prod_{j=1}^d (1 + \tau_j)\right)^{-n} \left(\prod_{j=1}^D \sigma_j\right)^{-n} \right) - \\ &\sum_{j=1}^d \frac{1}{2\sigma_j^2} \left(s_{1j} + \frac{1}{\tau_j^2} s_{2j} \right) - \sum_{j=d+1}^D \frac{1}{2\sigma_j^2} \left(s_{1j} + s_{2j} \right). \end{aligned}$$

We fix m , W and maximize the log-likelihood function over τ and σ . In such a case we have to solve the following system of equations

$$\frac{\partial \ln(L(X; m, W, \sigma, \tau))}{\partial \sigma_j} = 0, \quad \frac{\partial \ln(L(X; m, W, \sigma, \tau))}{\partial \tau_j} = 0,$$

for $j = 1, \dots, D$. Hence

$$\begin{aligned} -\frac{n}{\sigma_j} + \sigma_j^{-3}(s_{1j} + \tau_j^{-2}s_{2j}) &= 0, \text{ for } j = 1, \dots, d, \\ -\frac{n}{\sigma_j} + \sigma_j^{-3}(s_{1j} + s_{2j}) &= 0, \text{ for } j > d, \\ -\frac{n}{1+\tau_j} + \frac{s_{2j}}{\tau_j^3 \sigma_j^2} &= 0, \text{ for } j = 1, \dots, d. \end{aligned}$$

By simple calculations we obtain the expressions for the estimators in 3. Substituting it into the log-likelihood function, we get $\hat{L}(m, W) =$

$$\begin{aligned} &= \left(\frac{2}{\pi}\right)^{\frac{dn}{2}} \left(\frac{1}{2\pi}\right)^{\frac{(D-d)n}{2}} |\det(W)|^n \left(\prod_{j=1}^d \frac{1}{\sqrt{n}} g_j(m, W)^{\frac{3}{2}}\right)^{-n} \\ &e^{-\frac{dn}{2}} \left(\prod_{j=d+1}^D \left(\frac{s_{1j}+s_{2j}}{n}\right)^{\frac{1}{2}}\right)^{-n} = \frac{2^{(d-D/2)n} n^{dn/2}}{(\pi e)^{Dn/2}}. \\ &\left(\frac{1}{|\det(W)|^{\frac{2}{3}}} \prod_{j=1}^d g_j(m, W)\right)^{-\frac{3n}{2}} \left(\prod_{j=d+1}^D \frac{(s_{1j}+s_{2j})}{n}\right)^{-\frac{n}{2}} \end{aligned}$$

□

6 Appendix B

We will need the following well-known lemma (for the convenience of the reader we provide the proof).

Lemma 6.1. *Let $A = (a_{ij})_{1 \leq i, j \leq d}$ be a differentiable map from real numbers to $d \times d$ matrices then*

$$\frac{\partial \det(A)}{\partial a_{ij}} = \text{adj}^T(A)_{ij}, \quad (5)$$

where $\text{adj}(A)$ 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

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

□

Proof of Theorem 3.3. Let us start with the partial derivative of $\ln(l)$ with respect to m . We have

$$\begin{aligned} \frac{\partial \ln l(X; m, W)}{\partial m_k} &= \sum_{j=1}^d \frac{\partial \ln(g_j(m, W))}{\partial m_k} + \sum_{j=d+1}^D \frac{\partial \ln((s_{1j}+s_{2j})^{\frac{1}{3}})}{\partial m_k} \\ &= \sum_{j=1}^d \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \frac{\partial (s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}})}{\partial m_k} + \sum_{j=d+1}^D \frac{1}{(s_{1j}+s_{2j})^{\frac{1}{3}}} \frac{\partial ((s_{1j}+s_{2j})^{\frac{1}{3}})}{\partial m_k} \\ &= \sum_{j=1}^d \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \left(\frac{1}{3s_{1j}^{\frac{2}{3}}} \frac{\partial s_{1j}}{\partial m_k} + \frac{1}{3s_{2j}^{\frac{2}{3}}} \frac{\partial s_{2j}}{\partial m_k} \right) \\ &+ \sum_{j=d+1}^D \frac{1}{(s_{1j}+s_{2j})^{\frac{1}{3}}} \frac{1}{3(s_{1j}+s_{2j})^{\frac{2}{3}}} \left(\frac{\partial s_{1j}}{\partial m_k} + \frac{\partial s_{2j}}{\partial m_k} \right). \end{aligned}$$

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

$$\begin{aligned}\frac{\partial s_{1j}}{\partial m_k} &= \sum_{i \in I_j} \frac{\partial [\omega_j^T(x_i - m)]^2}{\partial m_k} = \\ \sum_{i \in I_j} 2\omega_j^T(x_i - m) \frac{\partial \omega_j^T(x_i - m)}{\partial m_k} &= \sum_{i \in I_j} -2\omega_j^T(x_i - m)\omega_{jk}.\end{aligned}$$

Analogously we get

$$\frac{\partial s_{2j}}{\partial m_k} = \sum_{i \in I_j^c} -2\omega_j^T(x_i - m)\omega_{jk}.$$

Hence

$$\begin{aligned}\frac{\partial \ln l}{\partial m_k} &= \sum_{j=1}^d \frac{-1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \left(\frac{1}{3s_{1j}^{\frac{2}{3}}} \sum_{i \in I_j} 2\omega_j^T(x_i - m)\omega_{jk} + \right. \\ &\quad \left. \frac{1}{3s_{2j}^{\frac{2}{3}}} \sum_{i \in I_j^c} 2\omega_j^T(x_i - m)\omega_{jk} \right) + \sum_{j=d+1}^D \frac{-1}{3(s_{1j} + s_{2j})^{\frac{1}{3}}}. \\ &\quad \left(\sum_{i \in I_j} 2\omega_j^T(x_i - m)\omega_{jk} + \sum_{i \in I_j^c} 2\omega_j^T(x_i - m)\omega_{jk} \right).\end{aligned}$$

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

$$\frac{\partial \ln |\det(W)|^{-\frac{2}{3}}}{\partial \omega_{pk}} + \sum_{j=1}^d \frac{\partial \ln(g_j(m, W))}{\partial \omega_{pk}} + \sum_{j=d+1}^D \frac{\partial \ln((s_{1j} + s_{2j})^{\frac{1}{3}})}{\partial \omega_{pk}}.$$

To calculate the derivative of the determinant we use Jacobi's formula (see Lemma 6.1). Hence

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

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

$$\begin{aligned}\frac{\partial \ln(g_j(m, W))}{\partial \omega_{pk}} &= \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \frac{\partial (s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}})}{\partial \omega_{pk}} = \\ &\quad \frac{1}{s_{1j}^{\frac{1}{3}} + s_{2j}^{\frac{1}{3}}} \left(\frac{1}{3s_{1j}^{\frac{2}{3}}} \frac{\partial s_{1j}}{\partial \omega_{pk}} + \frac{1}{3s_{2j}^{\frac{2}{3}}} \frac{\partial s_{2j}}{\partial \omega_{pk}} \right),\end{aligned}$$

where

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

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

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

Moreover,

$$\begin{aligned}\frac{\partial \ln((s_{1j} + s_{2j})^{\frac{1}{3}})}{\partial \omega_{pk}} &= \frac{1}{(s_{1j} + s_{2j})^{\frac{1}{3}}} \frac{\partial ((s_{1j} + s_{2j})^{\frac{1}{3}})}{\partial \omega_{pk}} = \\ &\quad \frac{1}{(s_{1j} + s_{2j})^{\frac{1}{3}}} \frac{1}{3} \frac{1}{(s_{1j} + s_{2j})^{\frac{2}{3}}} \left(\frac{\partial s_{1j}}{\partial \omega_{pk}} + \frac{\partial s_{2j}}{\partial \omega_{pk}} \right),\end{aligned}$$

Hence we obtain

$$\begin{aligned} \frac{\partial \ln l}{\partial \omega_{pk}} = & -\frac{2}{3}(\omega^{-1})_{pk}^T + \\ & \frac{1}{s_{1p}^{\frac{1}{3}} + s_{2p}^{\frac{1}{3}}} \left(\frac{1}{3} s_{1p}^{-\frac{2}{3}} \sum_{i \in I_p} 2\omega_p^T(x_i - m)(x_{ik} - m_k) \right. \\ & \left. + \frac{1}{3} s_{2p}^{-\frac{2}{3}} \sum_{i \in I_p^c} 2\omega_p^T(x_i - m)(x_{ik} - m_k) \right) + \\ & \frac{1}{3(s_{1p} + s_{2p})} \left(\sum_{i \in I_p} 2\omega_p^T(x_i - m)(x_{ik} - m_k) + \right. \\ & \left. \sum_{i \in I_p^c} 2\omega_p^T(x_i - m)(x_{ik} - m_k) \right). \end{aligned}$$

□

References

- [1] Stéphanie Allasonniere and Laurent Younes. A stochastic algorithm for probabilistic independent component analysis. *The Annals of Applied Statistics*, pages 125–160, 2012.
- [2] H Attias and CE Schreiner. Blind source separation and deconvolution by dynamic component analysis. In *Neural Networks for Signal Processing [1997] VII. Proceedings of the 1997 IEEE Workshop*, pages 456–465. IEEE, 1997.
- [3] Christian F Beckmann. Modelling with independent components. *Neuroimage*, 62(2):891–901, 2012.
- [4] Christian F Beckmann and Stephen M Smith. Probabilistic independent component analysis for functional magnetic resonance imaging. *IEEE transactions on medical imaging*, 23(2):137–152, 2004.
- [5] Anthony J Bell and Terrence J Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural computation*, 7(6):1129–1159, 1995.
- [6] Adel Belouchrani, Jean-François Cardoso, et al. Maximum likelihood source separation by the expectation-maximization technique: Deterministic and stochastic implementation. In *Proc. Nolta*, volume 95, pages 49–53. Citeseer, 1995.
- [7] Gilles Blanchard, Motoaki Kawanabe, Masashi Sugiyama, Vladimir Spokoiny, and Klaus-Robert Mazller. In search of non-gaussian components of a high-dimensional distribution. *Journal of Machine Learning Research*, 7(Feb):247–282, 2006.
- [8] Cesar F Caiafa, Emanuele Salerno, Araceli N Proto, and L Fiumi. Blind spectral unmixing by local maximization of non-gaussianity. *Signal Processing*, 88(1):50–68, 2008.
- [9] Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley and Sons, 2012.
- [10] Christopher G Green, Rajesh R Nandy, and Dietmar Cordes. Pca-preprocessing of fmri data adversely affects the results of ica. In *Proceedings of international society of magnetic resonance in medicine*, volume 10, 2002.
- [11] Trevor Hastie and Rob Tibshirani. *ProDenICA: Product Density Estimation for ICA using tilted Gaussian density estimates*, 2010. R package version 1.0.

- [12] Nathaniel E. Helwig. *ica: Independent Component Analysis*, 2015. R package version 1.0-1.
- [13] Nathaniel E Helwig and Sungjin Hong. A critique of tensor probabilistic independent component analysis: implications and recommendations for multi-subject fmri data analysis. *Journal of neuroscience methods*, 213(2):263–273, 2013.
- [14] Aapo Hyvärinen. Fast and robust fixed-point algorithms for independent component analysis. *Neural Networks, IEEE Transactions on*, 10(3):626–634, 1999.
- [15] Aapo Hyvärinen, Juha Karhunen, and Erkki Oja. *Independent component analysis*, volume 46. John Wiley & Sons, 2004.
- [16] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4):411–430, 2000.
- [17] Sreeba John. The three-parameter two-piece normal family of distributions and its fitting. *Communications in Statistics-Theory and Methods*, 11(8):879–885, 1982.
- [18] J. Karvanen. *PearsonICA*, 2008. R package version 1.2-3.
- [19] Motoaki Kawanabe, Masashi Sugiyama, Gilles Blanchard, and Klaus-Robert Müller. A new algorithm of non-gaussian component analysis with radial kernel functions. *Annals of the Institute of Statistical Mathematics*, 59(1):57–75, 2007.
- [20] Urbano Lorenzo-Seva and Jos MF Ten Berge. Tucker’s congruence coefficient as a meaningful index of factor similarity. *Methodology*, 2(2):57–64, 2006.
- [21] Markus Matilainen, Jari Miettinen, Klaus Nordhausen, Hannu Oja, and Sara Taskinen. *tsBSS: Tools for Blind Source Separation for Time Series*, 2016. R package version 0.2.
- [22] Jari Miettinen, Klaus Nordhausen, Hannu Oja, and Sara Taskinen. Deflation-based fastica with adaptive choices of nonlinearities. *IEEE Transactions on Signal Processing*, 62(21):5716–5724, 2014.
- [23] Eric Moulines, J-F Cardoso, and Elisabeth Gassiat. Maximum likelihood for blind separation and deconvolution of noisy signals using mixture models. In *Acoustics, Speech, and Signal Processing, 1997. ICASSP-97., 1997 IEEE International Conference on*, volume 5, pages 3617–3620. IEEE, 1997.
- [24] Dinh Tuan Pham and Philippe Garat. Blind separation of mixture of independent sources through a quasi-maximum likelihood approach. *Signal Processing, IEEE Transactions on*, 45(7):1712–1725, 1997.
- [25] Benjamin B Risk, David S Matteson, and David Ruppert. Likelihood component analysis. *arXiv preprint arXiv:1511.01609*, 2015.
- [26] Alexander Samarov, Alexandre Tsybakov, et al. Nonparametric independent component analysis. *Bernoulli*, 10(4):565–582, 2004.
- [27] Richard J Samworth, Ming Yuan, et al. Independent component analysis via nonparametric maximum likelihood estimation. *The Annals of Statistics*, 40(6):2973–3002, 2012.
- [28] Ran Shi, Ying Guo, et al. Investigating differences in brain functional networks using hierarchical covariate-adjusted independent component analysis. *The Annals of Applied Statistics*, 10(4):1930–1957, 2017.

- [29] P. Spurek et al. Ica based on the data asymmetry. *Pattern Recognition (accepted)*, 2017.
- [30] Michael E Tipping and Christopher M Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622, 1999.
- [31] Mattias Villani and Rolf Larsson. The multivariate split normal distribution and asymmetric principal components analysis. *Communications in Statistics—Theory and Methods*, 35(6):1123–1140, 2006.
- [32] Nan Wang, Bo Du, Liangpei Zhang, and Lifu Zhang. An abundance characteristic-based independent component analysis for hyperspectral unmixing. *IEEE Transactions on Geoscience and Remote Sensing*, 53(1):416–428, 2015.
- [33] Feiyun Zhu, Ying Wang, Shiming Xiang and Bin Fan, and Chunhong Pan. Structured sparse method for hyperspectral unmixing. *ISPRS Journal of Photogrammetry and Remote Sensing*, 88:101–118, 2014.
- [34] Feiyun Zhu, Ying Wang, Bin Fan, Gaofeng Meng, and Chunhong Pan. Effective spectral unmixing via robust representation and learning-based sparsity. *CoRR*, abs/1409.0685, 2014.
- [35] Feiyun Zhu, Ying Wang, Bin Fan, Gaofeng Meng, Shiming Xiang, and Chunhong Pan. Spectral unmixing via data-guided sparsity. *CoRR*, abs/1403.3155, 2014.