

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/287249560>

# The squared symmetric FastICA estimator

Article in *Signal Processing* · February 2017

DOI: 10.1016/j.sigpro.2016.08.028 · Source: arXiv

CITATIONS

6

READS

54

5 authors, including:



**Jari Miettinen**

University of Turku

14 PUBLICATIONS 130 CITATIONS

[SEE PROFILE](#)



**Klaus Nordhausen**

TU Wien

83 PUBLICATIONS 650 CITATIONS

[SEE PROFILE](#)



**Hannu Oja**

University of Turku

245 PUBLICATIONS 4,741 CITATIONS

[SEE PROFILE](#)



**Sara Taskinen**

University of Jyväskylä

49 PUBLICATIONS 1,307 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



gllvms in ecology [View project](#)



Supervised dimension reduction [View project](#)

All content following this page was uploaded by [Sara Taskinen](#) on 17 March 2016.

The user has requested enhancement of the downloaded file.

# The squared symmetric FastICA estimator

Jari Miettinen, Klaus Nordhausen, Hannu Oja, Sara Taskinen and Joni Virta

**Abstract**—In this paper we study the theoretical properties of the deflation-based FastICA method, the original symmetric FastICA method, and a modified symmetric FastICA method, here called the squared symmetric FastICA. This modification is obtained by replacing the absolute values in the FastICA objective function by their squares. In the deflation-based case this replacement has no effect on the estimate since the maximization problem stays the same. However, in the symmetric case a novel estimate with unknown properties is obtained. In the paper we review the classic deflation-based and symmetric FastICA approaches and contrast these with the new squared symmetric version of FastICA. We find the estimating equations and derive the asymptotical properties of the squared symmetric FastICA estimator with an arbitrary choice of nonlinearity. Asymptotic variances of the unmixing matrix estimates are then used to compare their efficiencies for large sample sizes showing that the squared symmetric FastICA estimator outperforms the other two estimators in a wide variety of situations.

**Index Terms**—Affine equivariance, independent component analysis, limiting normality, minimum distance index

## I. INTRODUCTION

We assume that a  $p$ -variate random vector  $\mathbf{x} = (x_1, \dots, x_p)^T$  follows the basic independent component (IC) model, that is, the components of  $\mathbf{x}$  are linear mixtures of  $p$  mutually independent latent variables in  $\mathbf{z} = (z_1, \dots, z_p)^T$ . The model can then be written as

$$\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Omega}\mathbf{z}, \quad (1)$$

where  $\boldsymbol{\mu}$  is a location shift and  $\boldsymbol{\Omega}$  is a full-rank  $p \times p$  mixing matrix. In independent component analysis (ICA), parameter  $\boldsymbol{\mu}$  is usually regarded as a nuisance parameter as the main interest is to find, using a random sample  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  from the distribution of  $\mathbf{x}$ , an estimate for an unmixing matrix  $\boldsymbol{\Gamma}$  such that  $\boldsymbol{\Gamma}\mathbf{x}$  has independent components [7], [2], [3]. Note that versions of (1) also exist where the dimension of  $\mathbf{z}$  is larger than that of  $\mathbf{x}$  (the *underdetermined case*) or the other way around (the *overdetermined case*), in the latter of which we can simply apply a dimension reduction method at first stage. In this paper we, however, restrict to the case where  $\mathbf{x}$  and  $\mathbf{z}$  are of the same dimension.

The IC model (1) is a semiparametric model in the sense that the marginal distributions of the components  $z_1, \dots, z_p$  are unspecified. However, some assumptions on  $\mathbf{z}$  are needed in order to fix the model: For identifiability of  $\boldsymbol{\Omega}$ , we need to assume that

(A1) at most one of the components  $z_1, \dots, z_p$  is gaussian [18].

Nevertheless,  $\boldsymbol{\mu}$ ,  $\boldsymbol{\Omega}$  and  $\mathbf{z}$  are still confounded and the mixing matrix  $\boldsymbol{\Omega}$  can be identified only up to the order, the signs, and heterogenous multiplication of its columns. To fix  $\boldsymbol{\mu}$  and the scales of the columns of  $\boldsymbol{\Omega}$  we further assume that

(A2)  $E(z_i) = 0$  and  $E(z_i^2) = 1$  for  $i = 1, \dots, p$ .

After these assumptions, the order and signs of the columns of  $\boldsymbol{\Omega}$  still remain unidentified. For practical data analysis, this is, however, often sufficient. The impact of the component order on asymptotics is further discussed in Section III.

The solutions to the ICA problem are often formulated as algorithms with two steps. The first step is to whiten the data, and the second step is to find an orthogonal matrix that rotates the whitened data to independent components. In the following we formulate such an algorithm at the population level using the random variable  $\mathbf{x}$ : Let  $S(F_{\mathbf{x}}) = \text{Cov}(\mathbf{x})$  denote the covariance matrix of a random vector  $\mathbf{x}$ , where  $F_{\mathbf{x}}$  denotes the cumulative distribution function  $\mathbf{x}$ , and write  $\mathbf{x}_{st} = S^{-1/2}(F_{\mathbf{x}})(\mathbf{x} - E(\mathbf{x}))$  for the standardized (whitened) random vector. Here the square root matrix  $S^{-1/2}$  is chosen to be symmetric. The aim of the second step is to find the rows of an orthogonal matrix  $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_p)^T$ , either one by one (deflation-based approach) or simultaneously (symmetric approach). The symmetric version of the famous FastICA algorithm [6] finds the orthogonal matrix  $\mathbf{U}$ , which maximizes a measure of non-Gaussianity for the rotated components,

$$\sum_{j=1}^p |E[G(\mathbf{u}_j^T \mathbf{x}_{st})]|,$$

where  $G$  is a twice continuously differentiable, nonlinear and nonquadratic function (see Section II-E for more details).

In this paper we replace the absolute values by their squares and consider the objective function

$$\sum_{j=1}^p (E[G(\mathbf{u}_j^T \mathbf{x}_{st})])^2,$$

as suggested in [19], where the squared symmetric FastICA estimates based on convex combinations of the third and fourth squared cumulants were studied in detail. Notice that replacing the absolute values by their squares in the objective functions has been mentioned in [6] and [3, Section 6], but the idea was never carried further. In Section II we formulate unmixing matrix functionals based on the two symmetric approaches and the deflation-based approach. Some statistical properties of the old estimators are recalled in Section III, and the corresponding results of squared symmetric FastICA are derived for the first time for general function  $G$ . The efficiencies of the three estimators are compared in Section IV using both asymptotic results and simulations.

J. Miettinen and S. Taskinen are with the Department of Mathematics and Statistics, University of Jyväskylä, Jyväskylä, FIN-40014, Finland (e-mail: jari.p.miettinen@jyu.fi).

K. Nordhausen, H. Oja and J. Virta are with the Department of Mathematics and Statistics, University of Turku, Turku, FIN-20014, Finland.

## II. FASTICA FUNCTIONALS

In this section we give formal definitions of three different, two old and one new, FastICA unmixing matrix functionals with corresponding estimating equations and algorithms for their computation. The formal definition of the squared symmetric FastICA functional is new. The conditions for function  $G$  that ensure the consistency of the estimates is also discussed.

### A. IC functionals

Let again  $F_x$  denote the cumulative distribution function of a random vector  $x$  obeying the IC model (1), and write  $\Gamma(F_x)$  for the value of an unmixing matrix functional at the distribution  $F_x$ . Due to the ambiguity in model (1), it is natural to require that the separation result  $\Gamma(F_x)x = \Gamma(F_z)z$  does not depend on  $\mu$  and  $\Omega$  and the choice of  $z$  in the model specification. This is formalized in the following.

**Definition 1.** The  $p \times p$  matrix-valued functional  $\Gamma(F_x)$  is said to be independent component (IC) functional if

- 1)  $\Gamma(F_x)x$  has independent components for all  $x$  in the IC model (1), and
- 2)  $\Gamma(F_x)$  is affine-equivariant in the sense that  $\Gamma(F_{Ax+b}) = \Gamma(F_x)A^{-1}$  for all nonsingular  $p \times p$  full-rank matrices  $A$ , for all  $p$ -vectors  $b$  and for all  $x$  (even beyond the IC model).

The condition  $\Gamma(F_{Ax+b}) = \Gamma(F_x)A^{-1}$  can be relaxed to be true only up to permutations and sign changes of their rows. The corresponding sample version  $\hat{\Gamma} = \Gamma(\mathbf{X})$  is obtained when the IC functional is applied to the empirical distribution function of  $\mathbf{X} = (x_1, \dots, x_n)$ . Naturally, the estimator is then also affine equivariant in the sense that  $\Gamma(\mathbf{A}\mathbf{X} + \mathbf{b}\mathbf{1}_n^T) = \Gamma(\mathbf{X})\mathbf{A}^{-1}$  for all nonsingular  $p \times p$  full-rank matrices  $\mathbf{A}$  and for all  $p$ -vectors  $\mathbf{b}$ .

The rest of this section focuses on three specific FastICA functionals. For recent overviews of FastICA and its variants see also [9] and [22].

### B. Deflation-based approach

Deflation-based FastICA functional is based on the algorithm proposed in [4] and [6]. In deflation-based FastICA method the rows of an unmixing matrix are extracted one after another. The method can thus be used in situations where only the few most important components are needed. The statistical properties of the deflation-based method were studied in [16] and [17], where the influence functions and limiting variances and covariances of the rows of unmixing matrix were derived.

Assume now that  $x$  is an observation from an IC model (1) with mean vector  $\mu = E(x)$  and covariance matrix  $\mathbf{S} = \text{Cov}(x)$ . In deflation-based FastICA, the unmixing matrix  $\Gamma = (\gamma_1, \dots, \gamma_p)^T$  is estimated so that after finding  $\gamma_1, \dots, \gamma_{j-1}$ , the  $j$ th row vector  $\gamma_j$  maximizes a measure of non-Gaussianity

$$|E[G(\gamma_j^T(x - E(x)))]|$$

under the constraints  $\gamma_l^T \mathbf{S} \gamma_j = \delta_{lj}$ ,  $l = 1, \dots, j$ , where  $\delta_{lj}$  is the Kronecker delta  $\delta_{lj} = 1$  (0) as  $l = j$  ( $l \neq j$ ). The

requirements for the function  $G$  and the conventional choices of it are discussed in Section II-E.

The deflation-based FastICA functional  $\Gamma^d$  satisfies the following  $p$  estimating equations [17], [15]:

**Definition 2.** The deflation-based FastICA functional  $\Gamma^d = (\gamma_1^d, \dots, \gamma_p^d)^T$  solves the estimating equations

$$\mathbf{T}(\gamma_j) = \mathbf{S} \left( \sum_{l=1}^j \gamma_l \gamma_l^T \right) \mathbf{T}(\gamma_j), \quad j = 1, \dots, p,$$

where

$$\mathbf{T}(\gamma) = E[g(\gamma^T(x - E(x)))(x - E(x))],$$

and  $g = G'$ .

The estimating equations imply that  $\Gamma \mathbf{S} \Gamma^T = \mathbf{I}_p$ , that is,  $\Gamma = \mathbf{U} \mathbf{S}^{-1/2}$  for some orthogonal matrix  $\mathbf{U}$ . The estimation problem can then be reduced to the estimation of the rows of  $\mathbf{U}$  one by one. This suggests the following fixed-point algorithm for  $u_j$ :

$$\begin{aligned} u_j &\leftarrow \mathbf{T}(u_j) \\ u_j &\leftarrow \left( \mathbf{I}_p - \sum_{l=1}^{j-1} u_l u_l^T \right) u_j \\ u_j &\leftarrow \|u_j\|^{-1} u_j, \end{aligned}$$

where  $\mathbf{T}(u) = E[g(u^T x_{st})x_{st}]$  and  $x_{st}$  is the whitened random variable. However, this algorithm is unstable and we recommend the use of the original algorithm [4], a modified Newton-Raphson algorithm, where the first step is

$$u_j \leftarrow E[g(u_j^T x_{st})x_{st}] - E[g'(u_j^T x_{st})]u_j.$$

For the estimate based on the observed data set, all the expectations above are replaced by the sample averages, e.g.,  $E(x)$  is replaced by  $\bar{x}$  and  $\mathbf{S}$  by the sample covariance matrix  $\hat{\mathbf{S}}$ .

Notice that neither the estimating equations nor the algorithm fixes the order in which the components are found and the order to some extent depends on the initial value in the algorithm. Since a change in the estimation order changes the unmixing matrix estimate more than just by permuting its rows, deflation-based FastICA is not affine equivariant if the initial value is chosen randomly. To find an estimate which globally maximizes the objective function at each stage, we propose the following strategy to choose the initial value for the algorithm:

- 1) Find a preliminary consistent estimator  $\Gamma_0$  of  $\Gamma$ .
- 2) Find a permutation matrix  $\mathbf{P}$  such that  $|E[G((\mathbf{P}\Gamma_0\mathbf{x})_1)]| \geq \dots \geq |E[G((\mathbf{P}\Gamma_0\mathbf{x})_p)]|$ .
- 3) The orthogonal initial value for  $\mathbf{U}$  is  $\mathbf{P}\Gamma_0\mathbf{S}^{1/2}$ .

The preliminary estimate in step 1 can be for example k-JADE estimate [10]. This algorithm, as well as all other FastICA algorithms mentioned in this paper, are implemented in R package fICA [11].

The extraction order of the components is highly important not only for the affine equivariance of the estimate, but also for its efficiency. In the deflationary approach, accurate estimation

of the first components can be shown to have a direct impact on accurate estimation of the last components as well. [15] discussed the extraction order and the estimation efficiency and introduced the so-called reloaded deflation-based FastICA, where the extraction order is based on the minimization of the sum of the asymptotic variances, see Section III. [13] discussed the estimate that uses different G-functions for different components. Different versions of the algorithm and their performance analysis are presented, for example, in [24], [23].

### C. Symmetric approach

In symmetric FastICA approach, the rows of  $\mathbf{\Gamma} = (\gamma_1, \dots, \gamma_p)^T$  are found simultaneously by maximizing

$$\sum_{j=1}^p |E[G(\gamma_j^T(x - E(x)))]|$$

under the constraint  $\mathbf{\Gamma}\mathbf{S}\mathbf{\Gamma}^T = \mathbf{I}_p$ . The unmixing matrix  $\mathbf{\Gamma}$  optimizes the Lagrangian function

$$L(\mathbf{\Gamma}, \mathbf{\Theta}) = \sum_{j=1}^p |E[G(\gamma_j^T(x - E(x)))]| - \sum_{j=1}^p \theta_{jj}(\gamma_j^T \mathbf{S} \gamma_j - 1) - \sum_{j=1}^{p-1} \sum_{l=j+1}^p \theta_{lj} \gamma_l^T \mathbf{S} \gamma_j,$$

where symmetric matrix  $\mathbf{\Theta} = [\theta_{lj}]$  contains  $p(p+1)/2$  Lagrangian multipliers. Differentiating the above function with respect to  $\gamma_j$  and setting the derivative to zero yields

$$E[g(\gamma_j^T(x - E(x)))(x - E(x))] s_j = 2\theta_{jj} \mathbf{S} \gamma_j + \sum_{l < j} \theta_{lj} \mathbf{S} \gamma_l + \sum_{l > j} \theta_{jl} \mathbf{S} \gamma_l,$$

where  $g = G'$  and  $s_j = \text{sign}(E[G(\gamma_j^T(x - E(x)))])$ . Then by multiplying both sides by  $\gamma_l^T$  we obtain  $\gamma_l^T E[g(\gamma_j^T(x - E(x)))(x - E(x))] s_j = \theta_{lj}$ , for  $l < j$ , and  $\gamma_l^T E[g(\gamma_j^T(x - E(x)))(x - E(x))] s_j = \theta_{jl}$ , for  $l > j$ . Hence the solution  $\mathbf{\Gamma}$  must satisfy the following estimating equations

**Definition 3.** The symmetric FastICA functional  $\mathbf{\Gamma}^s = (\gamma_1^s, \dots, \gamma_p^s)^T$  solves the estimating equations

$$\gamma_l^T \mathbf{T}(\gamma_j) s_j = \gamma_j^T \mathbf{T}(\gamma_l) s_l \quad \text{and} \quad \gamma_l^T \mathbf{S} \gamma_j = \delta_{lj},$$

where  $j, l = 1, \dots, p$ , and

$$\mathbf{T}(\gamma) = E[g(\gamma^T(x - E(x)))(x - E(x))],$$

$g = G'$ ,  $s_j = \text{sign}(E[G(\gamma_j^T(x - E(x)))])$  and  $\delta_{lj}$  is the Kronecker delta.

Again,  $\mathbf{\Gamma} = \mathbf{U}\mathbf{S}^{-1/2}$  for some orthogonal matrix  $\mathbf{U}$ . Then the estimation equations for  $\mathbf{U}$  are

$$\mathbf{u}_j^T \mathbf{T}(\mathbf{u}_j) s_j = \mathbf{u}_j^T \mathbf{T}(\mathbf{u}_l) s_l \quad \text{and} \quad \mathbf{u}_l^T \mathbf{u}_j = \delta_{lj},$$

where  $l, j = 1, \dots, p$ ,  $\mathbf{T}(\mathbf{u}) = E[g(\mathbf{u}^T \mathbf{x}_{st}) \mathbf{x}_{st}]$ , and the equations suggest the following fixed-point algorithm for  $\mathbf{U}$ :

$$\begin{aligned} \mathbf{T} &\leftarrow (\mathbf{T}(\mathbf{u}_1), \dots, \mathbf{T}(\mathbf{u}_p))^T \\ \mathbf{U} &\leftarrow (\mathbf{T}\mathbf{T}^T)^{-1/2} \mathbf{T}. \end{aligned}$$

As in the deflation-based approach, a more stable algorithm is obtained when  $\mathbf{T}(\mathbf{u}_j)$  is replaced by

$$\mathbf{T}^*(\mathbf{u}_j) = E[g(\mathbf{u}_j^T \mathbf{x}_{st}) \mathbf{x}_{st}] - E[g'(\mathbf{u}_j^T \mathbf{x}_{st})] \mathbf{u}_j.$$

In symmetric FastICA, different initial values give identical unmixing matrix estimates up to order and signs of the rows.

### D. Squared symmetric approach

In squared symmetric FastICA, the absolute values in the objective function of the regular symmetric FastICA are replaced by squares [19]. The squared symmetric FastICA functional  $\mathbf{\Gamma}^{s2} = (\gamma_1^{s2}, \dots, \gamma_p^{s2})^T$  maximizes

$$\sum_{j=1}^p (E[G(\gamma_j^T(x - E(x)))]^2)$$

under the constraint  $\mathbf{\Gamma}\mathbf{S}\mathbf{\Gamma}^T = \mathbf{I}_p$ . Similarly as in Section II-C the Lagrange multipliers method yields the following estimating equations:

**Definition 4.** The squared symmetric FastICA functional  $\mathbf{\Gamma}^{s2} = (\gamma_1^{s2}, \dots, \gamma_p^{s2})^T$  solves the estimating equations

$$\gamma_l^T \mathbf{T}_2(\gamma_j) = \gamma_j^T \mathbf{T}_2(\gamma_l) \quad \text{and} \quad \gamma_l^T \mathbf{S} \gamma_j = \delta_{lj},$$

where  $j, l = 1, \dots, p$ ,

$$\begin{aligned} \mathbf{T}_2(\gamma) &= E[G(\gamma^T(x - E(x)))] E[g(\gamma^T(x - E(x)))(x - E(x))], \\ g &= G' \text{ and } \delta_{lj} \text{ is the Kronecker delta.} \end{aligned}$$

The estimation equations for  $\mathbf{U}$  are

$$\mathbf{u}_l^T \mathbf{T}_2(\mathbf{u}_j) = \mathbf{u}_j^T \mathbf{T}_2(\mathbf{u}_l) \quad \text{and} \quad \mathbf{u}_l^T \mathbf{u}_j = \delta_{lj}, \quad l, j = 1, \dots, p,$$

where  $\mathbf{T}_2(\mathbf{u}) = E[G(\mathbf{u}^T \mathbf{x}_{st})] E[g(\mathbf{u}^T \mathbf{x}_{st}) \mathbf{x}_{st}]$ . The following algorithm, which is based on the same idea as the algorithm for symmetric FastICA, can be used to find the solution in practice:

$$\begin{aligned} \mathbf{T} &\leftarrow (\mathbf{T}_2^*(\mathbf{u}_1), \dots, \mathbf{T}_2^*(\mathbf{u}_p))^T \\ \mathbf{U} &\leftarrow (\mathbf{T}\mathbf{T}^T)^{-1/2} \mathbf{T}, \end{aligned}$$

where  $\mathbf{T}_2^*(\mathbf{u}) = E[G(\mathbf{u}^T \mathbf{x}_{st})] \{E[g(\mathbf{u}^T \mathbf{x}_{st}) \mathbf{x}_{st}] - E[g'(\mathbf{u}^T \mathbf{x}_{st})] \mathbf{u}\}$ .

Notice that

$$\mathbf{T}_2^*(\mathbf{u}) = E[G(\mathbf{u}^T \mathbf{x}_{st})] \mathbf{T}^*(\mathbf{u}),$$

and hence the squared symmetric FastICA estimator can be seen as weighted classical symmetric FastICA estimator. The more nongaussian, as measured by function  $G$ , an independent component is, the more impact it has in the orthogonalization step.

### E. Function $G$

The function  $G$  is required to be twice continuously differentiable, nonlinear and nonquadratic function such that  $E[G(z)] = 0$ , when  $z$  is a standard Gaussian random variable. The derivative function  $g = G'$  is the so-called nonlinearity. The use of classical kurtosis as a measure of non-Gaussianity is given by the nonlinearity function  $g(z) = z^3$  (pow3) [4].

Other popular choices include  $g(z) = \tanh(az)$  (*tanh*) and  $g(z) = z \exp(-az^2/2)$  (*gaus*) with tuning parameters  $a$  as suggested in [5], and  $g(z) = z^2$  (*skew*).

The deflation-based, symmetric and squared symmetric FastICA estimators need extra conditions for  $G$  to ensure the consistency of the estimation procedure: One then requires that, for any bivariate  $\mathbf{Z} = (z_1, z_2)^T$  with independent and standardized components ( $E(\mathbf{z}) = \mathbf{0}$  and  $\text{Cov}(\mathbf{z}) = \mathbf{I}_2$ ) and for any orthogonal  $2 \times 2$  matrix  $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2)^T$ ,

$$\begin{aligned} \text{def} \quad & |E[G(\mathbf{u}_1^T \mathbf{z})]| \leq \max(|E[G(z_1)]|, |E[G(z_2)]|), \\ \text{sym} \quad & |E[G(\mathbf{u}_1^T \mathbf{z})]| + |E[G(\mathbf{u}_2^T \mathbf{z})]| \\ & \leq |E[G(z_1)]| + |E[G(z_2)]| \\ \text{sym2} \quad & (E[G(\mathbf{u}_1^T \mathbf{z})])^2 + (E[G(\mathbf{u}_2^T \mathbf{z})])^2 \\ & \leq (E[G(z_1)])^2 + (E[G(z_2)])^2 \end{aligned}$$

[14] and [19] proved that for *pow3* and *skew* (as well as for their convex combination), all three conditions are satisfied. On the contrary, *tanh* and *gaus* do not satisfy the conditions for all choices of the distributions of  $z_1$  and  $z_2$ . For these two nonlinearities [20] found bimodal distributions for which the fixed points of the deflation-based FastICA algorithm are not correct solutions of the IC problem. In Figure 1 we plot the density functions of random variables  $z_1$  and  $z_2$  which serve as examples for a case where none of the three inequalities hold for *gaus*. These examples should however be seen as

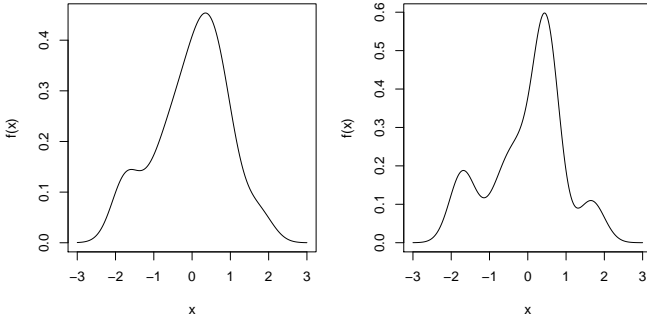


Fig. 1. Density functions of  $z_1$  and  $z_2$ , which violate the conditions def, sym and sym2 with nonlinearity *gaus*. Both distributions are mixtures of four Gaussian distributions. For more details, see Appendix.

rare and artificial exceptions and FastICA with *tanh* and *gaus* satisfy the conditions for most of pairs of distributions of  $\mathbf{z}$  we have checked. For example, in Section IV-C FastICA with *tanh* worked as expected under a wide variety of source distributions. Deflation-based or symmetric FastICA with *tanh* is perhaps the most popular unmixing matrix estimate.

See Section IV-B for the optimal choice of the nonlinearity for a component with a known density function.

### III. ASYMPTOTICAL PROPERTIES OF THE FASTICA ESTIMATORS

The limiting variances and the asymptotic multinormality of the deflation-based and symmetric FastICA unmixing matrix estimators were found quite recently in [17], [15], [21] and [22]. In this section, we review these findings and derive the results for the squared symmetric FastICA estimator.

Let now  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  be a random sample from the distribution of  $\mathbf{x}$  following the IC model (1). The deflation-based, symmetric and squared symmetric FastICA estimators  $\hat{\Gamma}^d$ ,  $\hat{\Gamma}^s$  and  $\hat{\Gamma}^{s2}$  are then obtained when the three functionals are applied to the empirical distribution of  $\mathbf{X}$ .

Due to affine equivariance, we can in the following assume without loss of generality that  $\Omega = \mathbf{I}_p$ . Before proceeding we need to make some additional assumptions on the distribution of  $\mathbf{z}_i = (z_{i1}, \dots, z_{ip})^T$ , namely,

(A3) The fourth moments  $\beta_j = E[z_{ij}^4]$  as well as the following expected values

$$\begin{aligned} \nu_j &= E[G(z_{ij})], & \mu_j &= E[g(z_{ij})], & \sigma_j^2 &= \text{Var}[g(z_{ij})], \\ \lambda_j &= E[g(z_{ij})z_{ij}], & \delta_j &= E[g'(z_{ij})], & \tau_j &= E[g'(z_{ij})z_{ij}] \end{aligned}$$

exist. Write also  $s_j = \text{sign}(\nu_j)$ .

Write now

$$\begin{aligned} \mathbf{T}_j &= \frac{1}{n} \sum_{i=1}^n (g(z_{ij}) - \mu_j) \mathbf{z}_i \quad \text{and} \\ \mathbf{T}_{2j} &= \frac{1}{n} \sum_{i=1}^n G(z_{ij}) \frac{1}{n} \sum_{i=1}^n (g(z_{ij}) - \mu_j) \mathbf{z}_i \end{aligned}$$

for  $j = 1, \dots, p$ . To avoid division by zero in the following theorem, assume that  $\nu_j(\lambda_j - \delta_j) \geq 0$  for all  $j = 1, \dots, p$ , with equality for at most one  $j$ , see [6], who stated that  $\nu_j(\lambda_j - \delta_j) > 0$  for most of the reasonable functions  $G$  and distributions of  $z_{ij}$ . For (*pow3*),  $\nu_j(\lambda_j - \delta_j) > 0$  for any distribution with  $E(z_{ij}^4) \neq 3$ . The limiting behavior of the deflation-based FastICA estimate was first given in [15]. The corresponding results of the symmetrical FastICA estimates are given in the following. The result (iii) is proved in the Appendix and the proof of (ii) is essentially similar to that. In the following theorem,  $\mathbf{e}_i$  is a  $p$ -vector with  $i$ th element one and others zero and  $o_P(1)$  replaces a random variable that converges in probability to zero as  $n$  goes to the infinity.

**Theorem 1.** Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  be a random sample from the IC model (1) satisfying the assumptions (A1)-(A3): If  $\Omega = \mathbf{I}_p$  then there exist a sequence of solutions  $\hat{\Gamma}^d$ ,  $\hat{\Gamma}^s$  and  $\hat{\Gamma}^{s2}$  converging to  $\mathbf{I}_p$  such that

(i) (deflation-based)

$$\begin{aligned} \sqrt{n} \hat{\gamma}_{jl}^d &= -\sqrt{n} \hat{\gamma}_{lj}^d - \sqrt{n} \hat{\mathbf{S}}_{jl} + o_P(1), \quad l < j, \\ \sqrt{n} (\hat{\gamma}_{jj}^d - 1) &= -\frac{1}{2} \sqrt{n} (\hat{\mathbf{S}}_{jj} - 1) + o_P(1), \quad l = j, \\ \sqrt{n} \hat{\gamma}_{jl}^d &= \frac{\mathbf{e}_l^T \sqrt{n} \mathbf{T}_j - \lambda_j \sqrt{n} \hat{\mathbf{S}}_{jl}}{\lambda_j - \delta_j} + o_P(1), \quad l > j, \end{aligned}$$

(ii) (symmetric)

$$\begin{aligned} \sqrt{n} (\hat{\gamma}_{jj}^s - 1) &= -\frac{1}{2} \sqrt{n} (\hat{\mathbf{S}}_{jj} - 1) + o_P(1), \quad l = j, \\ \sqrt{n} \hat{\gamma}_{jl}^s &= \frac{\mathbf{e}_l^T \sqrt{n} \mathbf{T}_j s_j - \mathbf{e}_j^T \sqrt{n} \mathbf{T}_l s_l - (\lambda_j s_j - \delta_l s_l) \sqrt{n} \hat{\mathbf{S}}_{jl}}{(\lambda_j - \delta_j) s_j + (\lambda_l - \delta_l) s_l} \\ &\quad + o_P(1), \quad l \neq j, \end{aligned}$$

(iii) (squared symmetric)

$$\sqrt{n}(\hat{\gamma}_{jj}^{s2} - 1) = -\frac{1}{2}\sqrt{n}(\hat{\mathbf{S}}_{jj} - 1) + o_P(1), \quad l = j,$$

$$\begin{aligned} \sqrt{n}\hat{\gamma}_{jl}^{s2} &= \frac{\mathbf{e}_l^T \sqrt{n}\mathbf{T}_{2j} - \mathbf{e}_j^T \sqrt{n}\mathbf{T}_{2l} + (\nu_l \delta_l - \nu_j \lambda_j) \sqrt{n}\hat{\mathbf{S}}_{jl}}{\nu_j(\lambda_j - \delta_j) + \nu_l(\lambda_l - \delta_l)} \\ &\quad + o_P(1), \quad l \neq j. \end{aligned}$$

For the asymptotical properties of deflation-based FastICA for several nonlinearities  $g$ , see [13]. As seen from Theorem 1 (i), the limiting distributions of vectors  $\hat{\gamma}_1^d, \dots, \hat{\gamma}_p^d$  depend on the order in which they are found. It is shown in Corollary 2 that, for  $j < l$ , the asymptotic variances of  $\hat{\gamma}_{lj}^d$  and  $\hat{\gamma}_{jl}^d$  are equal and depend only on the distribution of the  $j$ th independent component. The limiting distributions of the diagonal elements do not depend on the method or the chosen nonlinearity  $g$ . [19] discovered that the squared symmetric FastICA estimator with (*pow3*) nonlinearity has the same asymptotics as the JADE (joint approximate diagonalization of eigenmatrices) estimator [1]. We then have the following straightforward but important corollaries.

**Corollary 1.** *Under the assumptions of Theorem 1, if the joint limiting distribution of  $\sqrt{n}\mathbf{T}_{jl}$  and  $\sqrt{n}\mathbf{T}_{2jl}$  for  $j \neq l = 1, \dots, p$  and  $\sqrt{n}(\hat{\mathbf{S}}_{jl} - \delta_{jl})$  for  $j, l = 1, \dots, p$ , is a multivariate normal distribution, then also the limiting distributions of  $\sqrt{n} \text{vec}(\hat{\mathbf{\Gamma}}^d - \mathbf{I}_p)$ ,  $\sqrt{n} \text{vec}(\hat{\mathbf{\Gamma}}^s - \mathbf{I}_p)$  and  $\sqrt{n} \text{vec}(\hat{\mathbf{\Gamma}}^{s2} - \mathbf{I}_p)$  are multivariate normal.*

**Corollary 2.** *Under the assumptions of Theorem 1, the asymptotic covariance matrix (ASV) of the  $j$ th source vectors are given by*

$$\begin{aligned} \text{ASV}(\hat{\gamma}_j^d) &= \sum_{l=1}^p \text{ASV}(\hat{\gamma}_{jl}^d) \mathbf{e}_l \mathbf{e}_l^T, \\ \text{ASV}(\hat{\gamma}_j^s) &= \sum_{l=1}^p \text{ASV}(\hat{\gamma}_{jl}^s) \mathbf{e}_l \mathbf{e}_l^T, \quad \text{and} \\ \text{ASV}(\hat{\gamma}_j^{s2}) &= \sum_{l=1}^p \text{ASV}(\hat{\gamma}_{jl}^{s2}) \mathbf{e}_l \mathbf{e}_l^T, \end{aligned}$$

where

(i) (deflation-based)

$$\begin{aligned} \text{ASV}(\hat{\gamma}_{jl}^d) &= \frac{\sigma_l^2 - \lambda_l^2}{(\lambda_l - \delta_l)^2} + 1, \quad l < j \\ \text{ASV}(\hat{\gamma}_{jj}^d) &= \frac{\beta_j - 1}{4}, \quad l = j \\ \text{ASV}(\hat{\gamma}_{jl}^d) &= \frac{\sigma_j^2 - \lambda_j^2}{(\lambda_j - \delta_j)^2}, \quad l > j. \end{aligned}$$

(ii) (symmetric)

$$\begin{aligned} \text{ASV}(\hat{\gamma}_{jj}^s) &= \frac{\beta_j - 1}{4}, \quad l = j \\ \text{ASV}(\hat{\gamma}_{jl}^s) &= \frac{\sigma_j^2 + \sigma_l^2 - \lambda_j^2 + \delta_l(\delta_l - 2\lambda_l)}{((\lambda_j - \delta_j)s_j + (\lambda_l - \delta_l)s_l)^2}, \quad l \neq j. \end{aligned}$$

(iii) (squared symmetric)

$$\begin{aligned} \text{ASV}(\hat{\gamma}_{jj}^{s2}) &= \frac{\beta_j - 1}{4}, \quad l = j \\ \text{ASV}(\hat{\gamma}_{jl}^{s2}) &= \frac{\nu_j^2(\sigma_j^2 - \lambda_j^2) + \nu_l^2(\sigma_l^2 + \delta_l(\delta_l - 2\lambda_l))}{(\nu_j(\lambda_j - \delta_j) + \nu_l(\lambda_l - \delta_l))^2}, \\ &\quad l \neq j. \end{aligned}$$

The asymptotic variances of the deflation-based and symmetric FastICA estimators were first derived in [17] and [21], respectively. The asymptotic covariance matrices of the FastICA estimators for given marginal densities can be computed using the R package BSSasymp [12].

#### IV. EFFICIENCY COMPARISONS

The asymptotical results derived in Section III allow us to evaluate and compare the performances of the FastICA methods. In this section the asymptotic and finite sample efficiencies of deflation-based and symmetric FastICA estimators are compared to those of squared symmetric FastICA estimators using a wide range of distributions with varying skewness and kurtosis values.

##### A. Performance index

We measure the finite sample performance of the unmixing matrix estimates using the minimum distance index [8]

$$\hat{D} = D(\hat{\mathbf{\Gamma}}\mathbf{\Omega}) = \frac{1}{\sqrt{p-1}} \inf_{\mathcal{C} \in \mathcal{C}} \|\mathbf{C}\hat{\mathbf{\Gamma}}\mathbf{\Omega} - \mathbf{I}_p\| \quad (2)$$

where  $\|\cdot\|$  is the matrix (Frobenius) norm and  $\mathcal{C}$  is the set of  $p \times p$  matrices with exactly one non-zero element in each column and each row. The minimum distance index is scaled so that  $0 \leq \hat{D} \leq 1$ . If  $\mathbf{\Omega} = \mathbf{I}_p$  and  $\sqrt{n} \text{vec}(\hat{\mathbf{\Gamma}} - \mathbf{I}_p) \rightarrow N_{p^2}(\mathbf{0}, \mathbf{\Sigma})$ , then the limiting distribution of  $n(p-1)\hat{D}^2$  is that of weighted sum of independent chi squared variables with the expected value

$$\text{Trace}[(\mathbf{I}_{p^2} - \mathbf{D}_{p,p})\mathbf{\Sigma}(\mathbf{I}_{p^2} - \mathbf{D}_{p,p})], \quad (3)$$

where  $\mathbf{D}_{p,p} = \sum_i (\mathbf{e}_i \mathbf{e}_i^T) \otimes (\mathbf{e}_i \mathbf{e}_i^T)$ , and  $\otimes$  means the Kronecker product. Notice that (3) equals the sum of the limiting variances of the off-diagonal elements of  $\sqrt{n} \text{vec}(\hat{\mathbf{\Gamma}} - \mathbf{I}_p)$  and therefore

$$\sum_{j=1}^{p-1} \sum_{l=j+1}^p (\text{ASV}(\hat{\gamma}_{jl}) + \text{ASV}(\hat{\gamma}_{lj})) \quad (4)$$

provides a global measure of the variation of the estimate  $\hat{\mathbf{\Gamma}}$ .

##### B. Asymptotic efficiency

Let  $f_j$  be the density function and  $g_j = -f'_j/f_j$  be the optimal location score function for the  $j$ th independent component  $z_j$ . Also let  $I_j = \text{Var}(g_j(z_j))$  be the Fisher information number for the location problem. Write

$$\alpha_j := \frac{\sigma_j^2 - \lambda_j^2}{(\lambda_j - \delta_j)^2} = [(I_j - 1)\rho_{g(z_j)g_j(z_j) \cdot z_j}^2]^{-1},$$

where  $\rho_{g(z_j)g_j(z_j) \cdot z_j}^2$  is the squared partial correlation between  $g(z_j)$  and  $g_j(z_j)$  given  $z_j$ . Then we have the following.

**Theorem 2.** For our three estimates and for non-gaussian  $z_j$  and  $z_l$ ,  $j \neq l$ ,  $ASV(\hat{\gamma}_{jl}) + ASV(\hat{\gamma}_{lj})$  is

$$\left(\frac{\beta_j}{\beta_j + \beta_l}\right)^2 (2\alpha_j + 1) + \left(\frac{\beta_l}{\beta_j + \beta_l}\right)^2 (2\alpha_l + 1)$$

where

$$\begin{cases} \beta_j = 1, & s_j(\lambda_j - \delta_j), \text{ and } \nu_j(\lambda_j - \delta_j) \\ \beta_l = 0, & s_l(\lambda_l - \delta_l), \text{ and } \nu_l(\lambda_l - \delta_l) \end{cases}$$

for deflation-based, symmetric and squared symmetric FastICA estimates, respectively.

Notice first that the value of  $ASV(\hat{\gamma}_{jl}) + ASV(\hat{\gamma}_{lj})$  only depends on the  $j$ th and  $l$ th marginal distributions, which means we can restrict the comparison to bivariate distributions as the other components have no impact. If the  $j$ th and  $l$ th marginal distributions are the same, then the three values of  $ASV(\hat{\gamma}_{jl}) + ASV(\hat{\gamma}_{lj})$  are

$$(2\alpha_j + 1), \quad \frac{1}{2}(2\alpha_j + 1) \quad \text{and} \quad \frac{1}{2}(2\alpha_j + 1)$$

and these are minimized with the choice  $g = g_j$ . So, if  $z_1, \dots, z_p$  are identically distributed with the density function  $f$ , then the optimal choice for  $g$  is  $-f'/f$ .

If the  $l$ th component is Gaussian then,  $\lambda_l = \delta_l$ , and for the deflation-based and squared symmetric FastICA estimates,  $ASV(\hat{\gamma}_{jl}) + ASV(\hat{\gamma}_{lj}) = (2\alpha_j + 1)$  and for the symmetric FastICA estimate one gets

$$\begin{aligned} ASV(\hat{\gamma}_{jl}) + ASV(\hat{\gamma}_{lj}) &= (2\alpha_j + 1) + \frac{2(\sigma_l^2 - \lambda_l^2)}{\beta_j^2} \\ &= (2\alpha_j + 1) + \frac{2\sigma_l^2}{\beta_j^2} \left(1 - \rho_{g(z_{il})z_{il}}^2\right) \end{aligned}$$

where  $\rho_{g(z_{il})z_{il}}$  is the correlation between  $g(z_{il})$  and  $z_{il}$ . The symmetric FastICA is therefore always poorest in this case.

For further comparison of the estimators we use two families of source distributions, the standardized exponential power distribution family and the standardized gamma distribution family. The density function of standardized exponential power distribution with shape parameter  $\beta$  is

$$f(x) = \frac{\beta \exp\{-(|x|/\alpha)^\beta\}}{2\alpha \Gamma(1/\beta)},$$

where  $\beta > 0$ ,  $\alpha = (\Gamma(1/\beta)/\Gamma(3/\beta))^{1/2}$  and  $\Gamma$  is the gamma function. The distribution is symmetric for any  $\beta$ , and  $\beta = 2$  gives the normal (Gaussian) distribution,  $\beta = 1$  gives the heavy-tailed Laplace distribution and the density converges to the low-tailed uniform distribution as  $\beta \rightarrow \infty$ . The density function of standardized gamma distribution with shape parameter  $\alpha$  is

$$f(x) = \frac{(x + \sqrt{\alpha})^{\alpha-1} \alpha^{\alpha/2} \exp\{-(x + \sqrt{\alpha})\sqrt{\alpha}\}}{\Gamma(\alpha)}.$$

Gamma distributions are right skew, and for  $\alpha = k/2$ , the distribution is a chi square distribution with  $k$  degrees of freedom,  $k = 1, 2, \dots$ . When  $\alpha = 1$ , we have an exponential distribution, and the distribution converges to a normal distribution as  $\alpha \rightarrow \infty$ .

We next compare the asymptotic variances of the unmixing matrix estimates with the same nonlinearity and for  $\Omega = \mathbf{I}_p$ . For the comparison, write

$$ARE_{s2,d} = \frac{ASV(\hat{\gamma}_{jl}^d) + ASV(\hat{\gamma}_{lj}^d)}{ASV(\hat{\gamma}_{jl}^{s2}) + ASV(\hat{\gamma}_{lj}^{s2})},$$

for the asymptotic relative efficiency of the squared symmetric estimate with respect to the deflation based estimate, and similarly for  $ARE_{s2,s}$ . Notice that  $ARE_{s2,d}$  and  $ARE_{s2,s}$  depend on the two marginal distribution as well as on the chosen nonlinearity. We then plot the contour maps of the ARE's as functions of the shape parameters of the exponential power or gamma distributions with nonlinearities *pow3* and *tanh*. The equal efficiency is given by the ARE value 1 and can be found using the bar with contour thresholds on the right-hand side of the figures.

As seen in Figure 2, the squared symmetric FastICA estimator is in most cases more efficient than the deflation-based estimator. In Figure 3 we use  $ARE_{s2,s}$  similarly for the comparison between symmetric and squared symmetric FastICA. In the figures, the darker the point the higher relative efficiency. Notice that  $ARE_{s2,d} = 1$  if one of the components is Gaussian, and  $ARE_{s2,s} = 1$  if  $E(G(z_1)) = E(G(z_2))$  (e.g. if the two distributions are the same). Figure 3 shows that the areas where  $ASV(\hat{\gamma}_{jl}^s) > ASV(\hat{\gamma}_{jl}^{s2})$  and  $ASV(\hat{\gamma}_{jl}^s) < ASV(\hat{\gamma}_{jl}^{s2})$  are almost equally large, but the differences in favour of the squared symmetric estimator are larger. Also, they occur in cases where the separation of the components is difficult, and hence the efficiency is important there.

In Table I the values of  $ARE_{s2,s}$  and  $ARE_{s2,d}$  are displayed for different pairs of source distributions and for *pow3* in the upper triangle and for *tanh* in the lower triangle. Table I presents a sample of the values of Figure 2 and Figure 3 in a numerical form.

### C. Finite-sample efficiencies

We compare the finite-sample efficiencies of the estimates in a simulation study using the same two-dimensional settings with  $\Omega = \mathbf{I}_p$  as in the previous section. In each setting we consider the average of  $n(p-1)\hat{D}^2$  which has limiting expected value  $ASV(\hat{\gamma}_{jl}) + ASV(\hat{\gamma}_{lj})$ . Thus, the simulation study also illustrates how well the asymptotic results approximate the finite-sample variances. Let  $\hat{\Gamma}_i^{s2}$  and  $\hat{\Gamma}_i^s$ ,  $i = 1, \dots, M$ , be the estimates from  $M$  samples of size  $n$ . Then the finite sample asymptotic relative efficiency is estimated by

$$\widehat{ARE}_{s2,s} = \frac{\sum_{i=1}^M \{D(\hat{\Gamma}_i^{s2} \Omega)^2\}}{\sum_{i=1}^M \{D(\hat{\Gamma}_i^s \Omega)^2\}}.$$

In Table II, we list the estimated values of  $ARE_{s2,s}$  and  $ARE_{s2,d}$  for the same set of distributions as in Table I. For each setting,  $M = 10000$  samples of size  $n = 1000$  are generated. In most of the settings, the ratios of the averages are close to the corresponding asymptotical values. When both components are nearly Gaussian, a larger sample size than 1000 is required for  $\widehat{ARE}_{s2,s}$  and  $\widehat{ARE}_{s2,d}$  to converge to  $ARE_{s2,s}$  and  $ARE_{s2,d}$ , respectively. Also, if  $E[G(z_{ij})] \approx E[G(z_{il})]$ ,

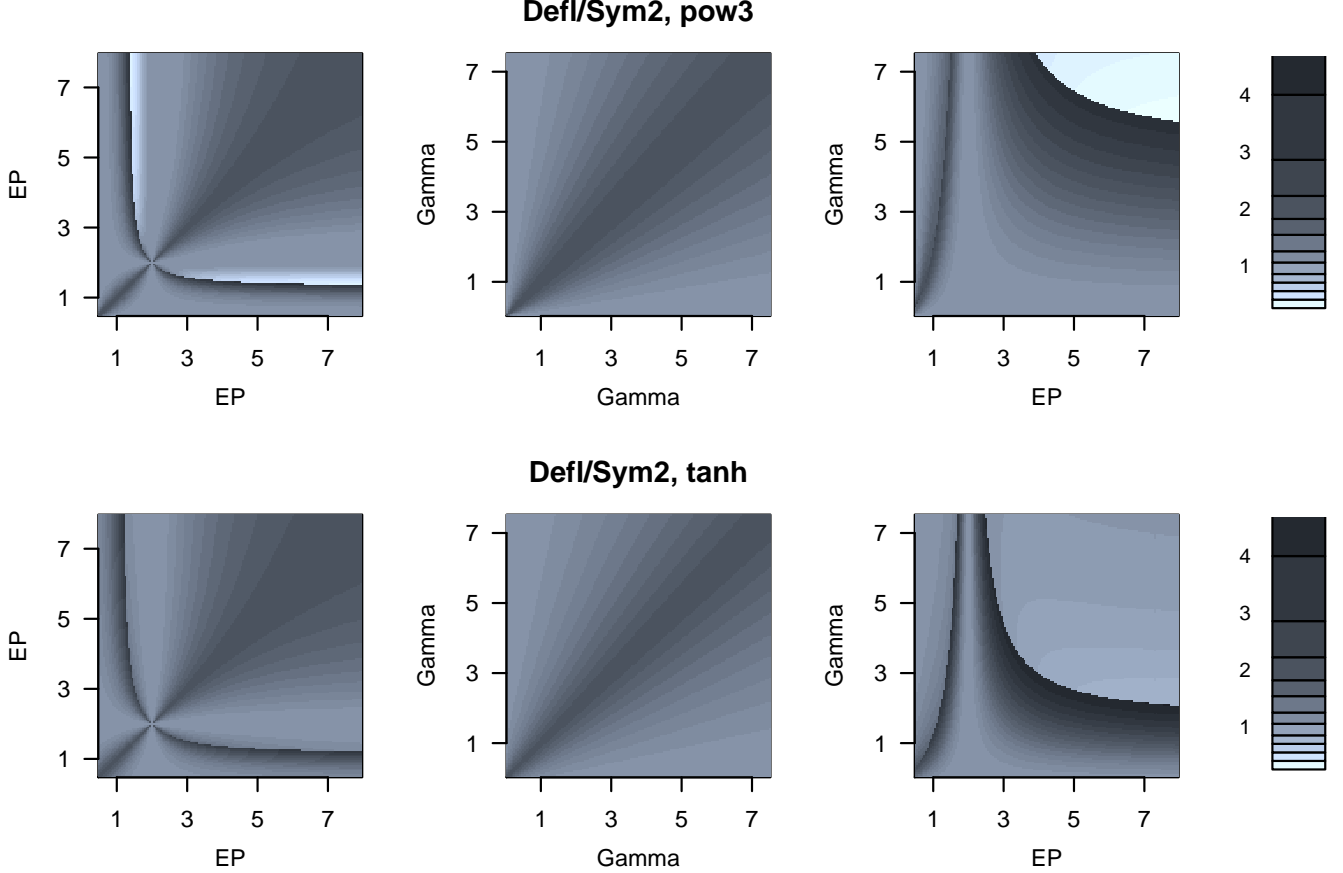


Fig. 2. Contour maps of asymptotic relative efficiencies  $ARE_{s2,d}$  when  $\Omega = \mathbf{I}_p$  and the source distributions are exponential power (EP) or gamma (Gamma) distributions with varying shape parameter values. The nonlinearities are *pow3* on the top row and *tanh* on the bottom row.

then the extraction order of the deflation-based estimate is not always the one which is assumed when computing the asymptotical variances. This may have a large impact on the efficiency of the deflation-based estimate.

In Figure 4 we plot the contour maps of the average of  $n(p-1)\hat{D}^2$  over 200 simulation runs for deflation-based, symmetric and squared symmetric FastICA estimates using *tanh*. Each setting has two independent components with exponential power distribution and varying shape parameter value, and  $n = 1000$ . Also the contour maps of the limiting expected values are given, and the corresponding maps resemble each other rather nicely. The asymptotical results thus provide good approximations already for  $n = 1000$ .

## V. CONCLUSIONS

In this paper we investigate in detail the properties of the squared symmetric FastICA procedure, obtained from the regular symmetric FastICA procedure by replacing in the objective function the analytically cumbersome absolute values by their squares. We reviewed in a unified way the estimating equations, algorithms and asymptotic theory of the classical deflation-based and symmetric FastICA estimators and provided similar tools and derived similar results for the novel squared symmetric FastICA. The asymptotic variances

were used to compare the three methods in numerous different situations.

The asymptotic and finite sample efficiency studies imply, that although none of the methods uniformly outperforms the others, the squared symmetric approach has the best overall performance under the considered combinations of source distributions and nonlinearities. Also, a crude ranking order of (*deflation-based*, *symmetric*, *squared symmetric*) from worst to best can be given and thus the use of the squared symmetric variant over the two other methods is highly recommended.

## APPENDIX

Write

$$\begin{aligned}\hat{T}(\hat{\gamma}) &= \frac{1}{n} \sum_{i=1}^n g(\hat{\gamma}^T(x_i - \bar{x}))(x_i - \bar{x}) \text{ and} \\ \hat{T}_2(\hat{\gamma}) &= \frac{1}{n} \sum_{i=1}^n G(\hat{\gamma}^T(x_i - \bar{x})) \frac{1}{n} \sum_{i=1}^n g(\hat{\gamma}^T(x_i - \bar{x}))(x_i - \bar{x}).\end{aligned}$$

The deflation-based, symmetric and squared symmetric FastICA estimators  $\hat{\Gamma}^d = \Gamma^d(\mathbf{X})$ ,  $\hat{\Gamma}^s = \Gamma^s(\mathbf{X})$  and  $\hat{\Gamma}^{s2} = \Gamma^{s2}(\mathbf{X})$  are defined as follows

**Definition 5.** The deflation-based FastICA estimate  $\hat{\Gamma}^d =$



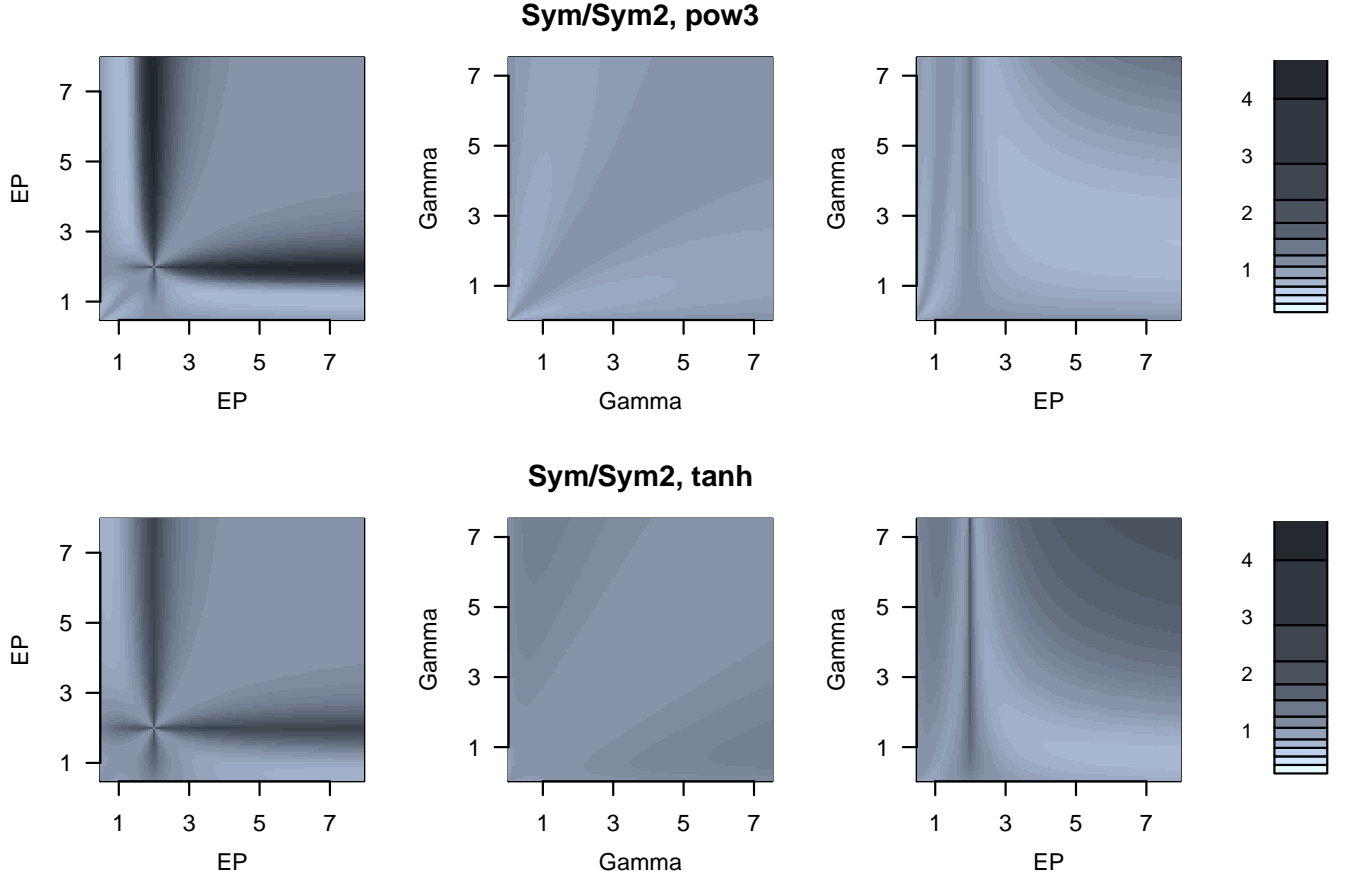


Fig. 3. Contour maps of asymptotic relative efficiencies  $ARE_{s2,s}$  when  $\Omega = \mathbf{I}_p$  and the source distributions are exponential power (EP) or gamma (Gamma) distributions with varying shape parameter values. The nonlinearities are *pow3* on the top row and *tanh* on the bottom row.

$(\hat{\gamma}_1^d, \dots, \hat{\gamma}_p^d)^T$  solves the estimating equations

$$\hat{\mathbf{T}}(\hat{\gamma}_j) = \hat{\mathbf{S}} \left( \sum_{l=1}^j \hat{\gamma}_l \hat{\gamma}_l^T \right) \hat{\mathbf{T}}(\hat{\gamma}_j), \quad j = 1, \dots, p, \quad (5)$$

**Definition 6.** The symmetric FastICA unmixing matrix estimate  $\hat{\mathbf{T}}^s = (\hat{\gamma}_1^s, \dots, \hat{\gamma}_p^s)^T$  solves the estimating equations

$$\hat{\gamma}_l^T \hat{\mathbf{T}}(\hat{\gamma}_j) \hat{s}_j = \hat{\gamma}_j^T \hat{\mathbf{T}}(\hat{\gamma}_l) \hat{s}_l \quad \text{and} \quad \hat{\gamma}_l^T \hat{\mathbf{S}} \hat{\gamma}_j = \delta_{lj}, \quad (6)$$

where  $j, l = 1, \dots, p$  and  $\delta_{lj}$  is the Kronecker delta.

**Definition 7.** The squared symmetric FastICA unmixing matrix estimate  $\hat{\mathbf{T}}^{s2} = (\hat{\gamma}_1^{s2}, \dots, \hat{\gamma}_p^{s2})^T$  solves the estimating equations

$$\hat{\gamma}_l^T \hat{\mathbf{T}}_2(\hat{\gamma}_j) = \hat{\gamma}_j^T \hat{\mathbf{T}}_2(\hat{\gamma}_l) \quad \text{and} \quad \hat{\gamma}_l^T \hat{\mathbf{S}} \hat{\gamma}_j = \delta_{lj}, \quad (7)$$

where  $j, l = 1, \dots, p$  and  $\delta_{lj}$  is the Kronecker delta.

To prove Theorem 1, we need the following straightforward result:

**Lemma 1.** The second set of estimating equations  $\hat{\gamma}_j^T \hat{\mathbf{S}} \hat{\gamma}_l = \delta_{lj}$ ,  $j, l = 1, \dots, p$  yields to

$$\sqrt{n}(\hat{\gamma}_{jj} - 1) = -\frac{1}{2}\sqrt{n}(\hat{\mathbf{S}} - \mathbf{I}_p)_{jj} + o_P(1)$$

and

$$\sqrt{n} \hat{\gamma}_{jl} + \sqrt{n} \hat{\gamma}_{lj} = -\sqrt{n} \hat{\mathbf{S}}_{jl} + o_P(1). \quad (8)$$

**Proof of Theorem 1 (iii)**

Let us now consider the first set of estimating equations. To shorten the notations, write  $\hat{\mathbf{T}}_2(\hat{\gamma}_j) = \hat{\mathbf{T}}_{2j}$ . Now

$$\sqrt{n} \hat{\gamma}_l^T \hat{\mathbf{T}}_{2j} = \sqrt{n} (\hat{\gamma}_l - \mathbf{e}_l)^T \hat{\mathbf{T}}_{2j} + \sqrt{n} \mathbf{e}_l^T (\hat{\mathbf{T}}_{2j} - \nu_j \lambda_j \mathbf{e}_j).$$

By Taylor expansion and Slutsky's Theorem, we have

$$\begin{aligned} \sqrt{n} (\hat{\mathbf{T}}_{2j} - \nu_j \lambda_j \mathbf{e}_j) &= \sqrt{n} (\mathbf{T}_{2j} - \nu_j \lambda_j \mathbf{e}_j) \\ &\quad - (\mu_j \lambda_j + \nu_j \tau_j) \mathbf{e}_j \mathbf{e}_j^T \sqrt{n} \bar{\mathbf{x}} \\ &\quad + (\lambda_j^2 \mathbf{e}_j \mathbf{e}_j^T + \nu_j \Delta_j) \sqrt{n} (\hat{\gamma}_j - \mathbf{e}_j) + o_P(1), \end{aligned}$$

where  $\Delta_j = E[g'(z_{ij}) \mathbf{z}_i \mathbf{z}_i^T]$ . Consequently,

$$\begin{aligned} \sqrt{n} \hat{\gamma}_l^T \hat{\mathbf{T}}_{2j} &= \sqrt{n} (\hat{\gamma}_l - \mathbf{e}_l)^T \nu_j \lambda_j \mathbf{e}_j \\ &\quad + \mathbf{e}_l^T (\sqrt{n} \mathbf{T}_{2j} - (\mu_j \lambda_j + \nu_j \tau_j) \mathbf{e}_j \mathbf{e}_j^T \sqrt{n} \bar{\mathbf{x}} \\ &\quad + (\lambda_j^2 \mathbf{e}_j \mathbf{e}_j^T + \nu_j \Delta_j) \sqrt{n} (\hat{\gamma}_j - \mathbf{e}_j)) + o_P(1) \\ &= \nu_j \lambda_j \sqrt{n} \hat{\gamma}_{lj} + \mathbf{e}_l^T \sqrt{n} \mathbf{T}_{2j} + \nu_j \delta_j \sqrt{n} \hat{\gamma}_{jl} + o_P(1). \end{aligned}$$

According to our estimating equations, above expression should be equivalent to

$$\sqrt{n} \hat{\gamma}_j^T \hat{\mathbf{T}}_{2l} = \nu_l \lambda_l \sqrt{n} \hat{\gamma}_{jl} + \sqrt{n} \mathbf{e}_j^T \mathbf{T}_{2l} + \nu_l \delta_l \sqrt{n} \hat{\gamma}_{lj} + o_P(1),$$

TABLE II

VALUES OF  $\widehat{ARE}_{s2,s}$  (ON THE TOP) AND  $\widehat{ARE}_{s2,d}$  (ON THE BOTTOM) COMPUTED FROM 10000 SAMPLES OF SIZE  $n = 1000$  FOR DIFFERENT DISTRIBUTIONS. L=LAPLACE=EP1=EXPONENTIAL POWER DISTRIBUTION WITH  $\beta = 1$ , N=NORMAL DISTRIBUTION=EP2, U=UNIFORM DISTRIBUTION, G1=GAMMA DISTRIBUTION WITH  $\alpha = 1$ . UPPER TRIANGLE FOR *pow3* AND LOWER TRIANGLE FOR *tanh*.

	L	EP1.5	EP1.75	N	EP 3	EP4	U	G1	G3	G6
L	0.98\0.86	0.94	1.09	1.17	0.80	0.75	0.73	0.85	0.88	0.98
EP1.5	1.03	1.01\0.95	1.15	1.34	0.94	1.07	1.49	0.92	0.89	0.95
EP1.75	1.27	1.27	1.14\1.13	1.07	1.66	2.10	3.13	1.02	1.06	1.08
N	1.47	1.65	1.12	–	2.22	2.99	4.13	1.05	1.25	1.25
EP3	0.92	1.06	1.59	1.91	1.31\1.07	1.08	1.45	0.85	0.77	0.90
EP4	0.84	1.14	1.67	2.12	1.06	1.07\1.00	1.14	0.81	0.76	0.98
U	0.87	1.41	1.89	2.22	1.25	1.08	1.00\1.00	0.76	0.82	1.23
G1	0.96	0.97	1.18	1.43	0.84	0.76	0.76	0.99\0.84	0.88	0.93
G3	1.15	0.94	1.12	1.51	0.87	1.01	1.38	1.13	1.10\0.83	0.90
G6	1.31	1.21	1.09	1.23	1.19	1.47	1.86	1.25	1.15	1.11\0.93
	L	EP1.5	EP1.75	N	EP 3	EP4	U	G1	G3	G6
L	1.93\1.79	1.12	1.02	1.00	1.09	1.18	1.38	1.47	1.53	1.18
EP1.5	1.14	1.82\1.47	1.14	0.98	1.66	1.54	0.85	1.04	1.29	1.40
EP1.75	1.03	1.24	1.14\1.04	1.03	1.18	0.89	0.72	1.00	1.04	1.11
N	1.00	0.98	1.05	–	0.89	0.91	0.92	1.00	1.00	1.00
EP3	1.18	1.78	1.19	0.93	1.82\1.78	1.35	1.01	1.03	1.25	1.50
EP4	1.42	1.41	1.01	0.97	1.40	1.93\1.94	1.23	1.05	1.42	1.52
U	2.14	1.00	0.99	0.99	1.13	1.33	1.92\1.94	1.11	1.64	1.30
G1	2.23	1.18	1.03	1.00	1.20	1.44	2.27	2.04\1.66	1.22	1.06
G3	1.43	1.99	1.20	0.99	1.79	1.85	1.03	1.26	2.18\1.57	1.36
G6	1.05	1.83	1.21	0.99	1.35	1.05	0.91	1.06	1.49	1.29\1.36

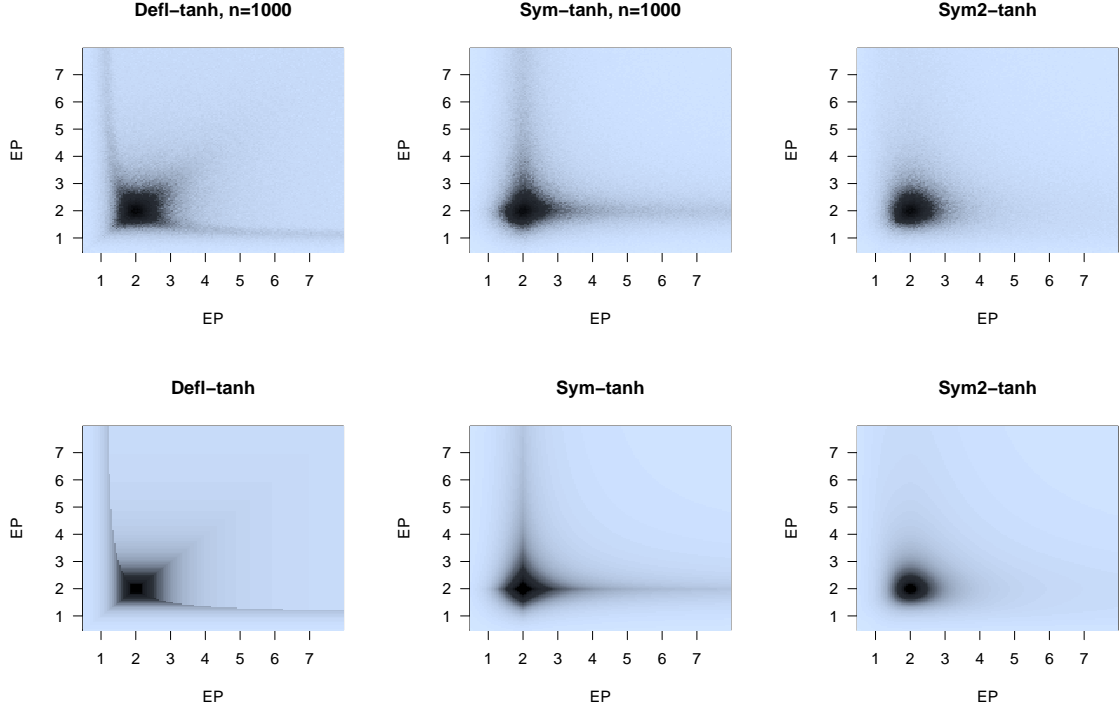


Fig. 4. Contour maps of the average of  $n(p-1)\hat{D}^2$  over 200 simulation runs with deflation-based, symmetric and squared symmetric FastICA estimates using *tanh* on the top and the contour maps of the limiting expected values on the bottom. Two independent components with exponential power distribution and varying shape parameter value.

which means that

$$\begin{aligned}
 & (\nu_l \lambda_l - \nu_j \delta_j) \sqrt{n} \hat{\gamma}_{jl} - (\nu_j \lambda_j - \nu_l \delta_l) \sqrt{n} \hat{\gamma}_{lj} \\
 &= \sqrt{n} \mathbf{e}_l^T \mathbf{T}_{2j} - \sqrt{n} \mathbf{e}_j^T \mathbf{T}_{2l} + o_P(1).
 \end{aligned}$$

Now using (8) in Lemma 1, we have that

$$\begin{aligned}
 & (\nu_l \lambda_l - \nu_j \delta_j) \sqrt{n} \hat{\gamma}_{jl} + (\nu_j \lambda_j - \nu_l \delta_l) \\
 & (\sqrt{n} \hat{\mathbf{S}}_{jl} + \sqrt{n} \hat{\gamma}_{jl}) = \sqrt{n} (\mathbf{e}_l^T \mathbf{T}_{2j} - \mathbf{e}_j^T \mathbf{T}_{2l}) + o_P(1).
 \end{aligned}$$

TABLE I

VALUES OF  $ARE_{s2,s}$  (ON THE TOP) AND  $ARE_{s2,d}$  (ON THE BOTTOM) FOR DIFFERENT DISTRIBUTIONS. L=LAPLACE=EP1=EXPONENTIAL POWER DISTRIBUTION WITH  $\beta = 1$ , N=NORMAL DISTRIBUTION=EP2, U=UNIFORM DISTRIBUTION, G1=GAMMA DISTRIBUTION WITH  $\alpha = 1$ . UPPER TRIANGLE FOR  $pow3$  AND LOWER TRIANGLE FOR  $tanh$ .

	L	EP1.5	EP1.75	N	EP 3	EP4	U	G1	G3	G6
L	1	0.87	0.96	1.10	0.79	0.74	0.70	0.84	1.00	0.94
EP1.5	1.05	1	1.02	1.46	0.87	1.05	1.61	0.87	0.84	0.95
EP1.75	1.21	1.10	1	1.72	1.53	2.21	3.75	0.95	0.93	0.92
N	1.50	1.79	1.90	–	3.09	3.85	5.49	1.03	1.13	1.25
EP3	0.91	0.99	1.29	2.23	1	1.23	1.45	0.86	0.74	0.77
EP4	0.85	1.13	1.55	2.32	1.05	1	1.13	0.81	0.71	0.86
U	0.87	1.46	1.90	2.45	1.24	1.08	1	0.76	0.74	1.20
G1	0.97	0.94	1.12	1.43	0.82	0.74	0.76	1	0.84	0.87
G3	1.15	0.93	0.97	1.67	0.85	1.18	1.39	1.07	1	0.93
G6	1.27	1.10	0.92	1.79	1.16	1.70	1.91	1.18	1.06	1
L	2	1.12	1.02	1	1.07	1.15	1.34	1.46	1.53	1.18
EP1.5	1.16	2	1.02	1	0.87	0.55	0.40	1.03	1.26	1.89
EP1.75	1.03	1.21	2	1	0.90	0.83	0.81	1.00	1.04	1.15
N	1	1	1	–	1	1	1	1	1	1
EP3	1.19	2.44	1.12	1	2	1.13	1.01	1.02	1.17	1.72
EP4	1.42	1.17	1.04	1	1.42	2	1.23	1.04	1.35	2.60
U	2.24	1.02	1.01	1	1.14	1.34	2	1.08	1.83	0.21
G1	2.32	1.18	1.03	1	1.19	1.74	2.24	2	1.20	1.05
G3	1.12	2.29	1.24	1	2.54	0.81	0.81	1.17	2	1.39
G6	1.04	1.15	1.91	1	0.94	0.93	0.94	1.05	1.29	2

Then

$$(\nu_j(\lambda_j - \delta_j) + \nu_l(\lambda_l - \delta_l))\sqrt{n}\hat{\gamma}_{jl} \\ = \sqrt{n}(\mathbf{e}_l^T \mathbf{T}_{2j} - \mathbf{e}_j^T \mathbf{T}_{2l}) + (\nu_l \delta_l - \nu_j \lambda_j)\sqrt{n}\hat{\mathbf{S}}_{jl} + o_P(1),$$

which proves the Theorem.

The densities of  $z_1$  and  $z_2$  in Section II-E are given by

$$f_i = \sum_{j=1}^4 \pi_{ij} N(\mu_{ij}, \sigma_{ij}^2), \quad i = 1, 2,$$

where  $N(\mu, \sigma^2)$  denotes the Gaussian density function with mean  $\mu$  and variance  $\sigma^2$ , and the (rounded) parameter values are

$$\begin{aligned} \pi_{11} &= 0.09, & \pi_{12} &= 0.43, & \pi_{13} &= 0.43, & \pi_{14} &= 0.04, \\ \pi_{21} &= 0.15, & \pi_{22} &= 0.31, & \pi_{23} &= 0.45, & \pi_{24} &= 0.09, \\ \mu_{11} &= -1.76, & \mu_{12} &= -0.34, & \mu_{13} &= 0.54, & \mu_{14} &= 1.79, \\ \mu_{21} &= -1.71, & \mu_{22} &= -0.36, & \mu_{23} &= 0.48, & \mu_{24} &= 1.66, \\ \sigma_{11}^2 &= 0.13, & \sigma_{12}^2 &= 0.50, & \sigma_{13}^2 &= 0.28, & \sigma_{14}^2 &= 0.13, \\ \sigma_{21}^2 &= 0.11, & \sigma_{22}^2 &= 0.26, & \sigma_{23}^2 &= 0.11, & \sigma_{24}^2 &= 0.11. \end{aligned}$$

#### ACKNOWLEDGMENT

This work was supported by the Academy of Finland (grants 251965, 256291 and 268703).

#### REFERENCES

- [1] J.C. Cardoso and A. Souloumiac, "Blind beamforming for non gaussian signals," *IEE Proceedings-F*, vol. 140, pp. 362–370, 1993.
- [2] A. Cichocki and S. Amari, *Adaptive Blind Signal and Image Processing: Learning Algorithms and Applications*. Wiley, Cichester, 2002.
- [3] P. Comon and C. Jutten, *Handbook of Blind Source Separation: Independent Component Analysis and Applications*. Academic Press, Oxford, 2010.
- [4] A. Hyvärinen and E. Oja, "A fast fixed-point algorithm for independent component analysis," *Neural Computation*, vol. 9, pp. 1483–1492, 1997.
- [5] A. Hyvärinen, "One-Unit Contrast Functions for Independent Component Analysis: A Statistical Analysis," in *Neural Networks for Signal Processing VII* (Proc. IEEE NNSP Workshop 1997), Amelia Island, Florida, pp. 388–397.
- [6] A. Hyvärinen, "Fast and Robust fixed-point algorithms for independent component analysis," *IEEE Trans. Neural Networks*, vol. 10, pp. 626–634, 1999.
- [7] A. Hyvärinen, J. Karhunen and E. Oja, *Independent Component Analysis*. John Wiley and Sons, New York, 2001.
- [8] P. Ilmonen, K. Nordhausen, H. Oja and E. Ollila, "A new performance index for ICA: properties computation and asymptotic analysis," in *Latent Variable Analysis and Signal Processing (Proceedings of 9th International Conference on Latent Variable Analysis and Signal Separation)*, 229–236, 2010.
- [9] Z. Koldovský and P. Tichavský, "Improved Variants of the FastICA Algorithm," in *Advances in Independent Component Analysis and Learning Machines*, Academic Press, Springer, 2015.
- [10] J. Miettinen, K. Nordhausen, H. Oja and S. Taskinen, "Fast Equivariant JADE," in *Proceedings of "38th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2013)"*, Vancouver, 2013, pp. 6153–6157.
- [11] J. Miettinen, K. Nordhausen, H. Oja and S. Taskinen, "fICA: Classical, Reloaded and Adaptive FastICA Algorithms," R package version 1.0-3, <http://cran.r-project.org/web/packages/fICA>, 2015.
- [12] J. Miettinen, K. Nordhausen, H. Oja and S. Taskinen, "BSSasymp: Asymptotic Covariance Matrices of Some BSS Mixing and Unmixing Matrix Estimates", R package version 1.1-1, <http://cran.r-project.org/web/packages/BSSasymp>, 2015.
- [13] J. Miettinen, K. Nordhausen, H. Oja and S. Taskinen, "Deflation-based FastICA with adaptive choices of nonlinearities", *IEEE Transactions on Signal Processing*, vol. 62, pp. 5716–5724, 2014.
- [14] J. Miettinen, S. Taskinen, K. Nordhausen and H. Oja, "Fourth moments and independent component analysis", *Statistical Science*, vol. 30, pp. 372–390, 2015.
- [15] K. Nordhausen, P. Ilmonen, A. Mandal, H. Oja and E. Ollila, "Deflation-based FastICA reloaded", in *Proc. "19th European Signal Processing Conference 2011 (EUSIPCO 2011)"*, Barcelona, 2011, pp. 1854–1858.
- [16] E. Ollila, "On the robustness of the deflation-based FastICA estimator", in *Proc. IEEE Workshop on Statistical Signal Processing (SSP'09)*, pp. 673–676, 2009.
- [17] E. Ollila, "The deflation-based FastICA estimator: statistical analysis revisited", *IEEE Transactions on Signal Processing*, vol. 58, pp. 1527–1541, 2010.
- [18] F.J. Theis, "A new concept for separability problems in blind source separation", *Neural Computation*, vol. 16, pp. 1827–1850, 2004.
- [19] J. Virta, K. Nordhausen and H. Oja, "Joint use of third and fourth cumulants in independent component analysis", arXiv preprint arXiv:1505.02613.
- [20] T. Wei, "On the spurious solutions of the Fastica algorithm", in *Statistical Signal Processing (SSP), 2014 IEEE Workshop on*, pp. 161–164, 2014.
- [21] T. Wei, "A convergence and asymptotic analysis of the generalized symmetric FastICA algorithm", *IEEE Transactions on Signal Processing*, vol. 63, pp. 6445–6458, 2015.
- [22] T. Wei, "An overview of the asymptotic performance of the family of the FastICA algorithms", in E. Vincent et al. (Eds.): *LVA/ICA 2015, LNCS 9237*, Springer, pp. 336–343, 2015.
- [23] V. Zarzoso and P. Comon, "Comparative Speed Analysis of FastICA", in M.E. Davies et al. (Eds.): *Independent Component Analysis and Signal Separation, LNCS 4666*, Springer, pp. 293–300, 2007.
- [24] V. Zarzoso, P. Comon, and M. Kallel M., "How fast is FastICA?" in *Proceedings for 14th European Signal Processing Conference*, pp. 1–5, 2006.